

Volume 68 • Number 4 • August 2020

# Acta Geophysica



PAN

INSTITUTE OF PHYSICS  
POLISH ACADEMY OF SCIENCES



Institute of Geophysics  
Polish Academy of Sciences



Springer



# Seismicity analysis of selected faults in Makran Southern Pakistan

Muhammad Jahangir Khan<sup>1</sup> · Mubarak Ali<sup>2</sup> · Min Xu<sup>3</sup> · Mehrab Khan<sup>1</sup>

Received: 13 January 2020 / Accepted: 15 May 2020 / Published online: 29 May 2020  
© Institute of Geophysics, Polish Academy of Sciences & Polish Academy of Sciences 2020

## Abstract

This study was focused on evaluation of seismicity yield of key faults guarding the terrestrial Makran region (Balochistan and Iran) using updated datasets. The Makran onshore region is bounded by paramount strike-slip faults (Chaman fault, Ghazaband fault, Ornach-Nal fault and Minab fault). We have compiled the earthquake's catalog and performed iterative processes for declustering of independent main shocks, estimation of magnitude completeness and  $b$ -values. The main shocks of disastrous earthquakes, e.g.,  $M_w$  7.7 Quetta (1935) and  $M_w$  7.7 Awaran (2013), have epicenters along the Ghazaband fault and splay of Chaman fault, respectively. The earthquake source parameters such as moment magnitude, focal depth, focal mechanism, and epicenter location was utilized in mapping and seismicity evaluation of the faults. The focal mechanism solution was derived to determine the fault mechanics in generating the  $M_w > 5.0$  events along these faults. This study helped us to compare the seismological profiles of each boundary fault and present the seismological information including the characteristics of events on/along fault, estimation of re-occurrence period, corresponding  $b$ -value, and hazard potential of the of the key faults. Since our study is based on recent dataset, i.e., inclusion of 2013  $M_w$  7.7 Awaran earthquake, the estimated results could help in better planning against the earthquake hazard in near-field cities & coastal towns of southern Pakistan.

**Keywords** Seismicity · GIS · Karachi · Southern Pakistan · Awaran earthquake

## Introduction

Recent studies enhance the seismogeological understanding about interseismic strain, seismogenic behavior and seismotectonic dynamics of Makran region, southern Pakistan (Nemati 2019; Burg 2018; Penney et al. 2017). The research studies (Khan 2015; Barnhart et al. 2014; Hadi et al. 2013; Smith et al. 2013; Quittmeyer and Kafka 1984) revealed that Makran region is an active seismotectonic zone; however, it is diversified in seismotectonic characteristics in its eastern and western parts such as strength of earthquakes, seismicity recurrence, frequency of events, clustering of after- and/or foreshocks, focal depths of main shocks from interplate and intraplate regions and corresponding  $b$ -values (Penney et al.

2017; Khan 2015; Ambraseys and Bilham 2014; Smith et al. 2013; Rani et al. 2011; Regard et al. 2010). Makran region distinguishes four main litho-tectonic units from north to south (offshore, costal, outer and inner Makran), and this geological subdivision is associated with strong deformation such as long wavelength folding (Burg 2018; McCall 2003) and secondary faulting (Dolati and Burg 2013) within massive basin fill deposits (~ 7 km, Smith 2013) of Makran. The tectonic lineaments and massive geological structures in Makran onshore are revealing the dynamic source of tectonic platform in southern Pakistan.

It was hypothesized that the major seismicity of the Makran onshore region (MKOR) is controlled by transform faults, i.e., Chaman fault (CF), Ghazaband fault (GF), Ornach-Nal fault (ONF), the Minab fault (MF) and their associated fault splays and megathrust of Makran subduction. MKOR is bounded by the paramount strike-slip/transform faults: Minab fault at western side and Chaman, Ghazaband, Ornach-Nal faults at eastern side (Fig. 1). We have focused on the seismicity of major strike-slip/transform faults of MKOR in this study. The aim of this study was to evaluate the seismological characteristics of these boundary faults using instrumental data (over ~ 40 years, 1978–2018).

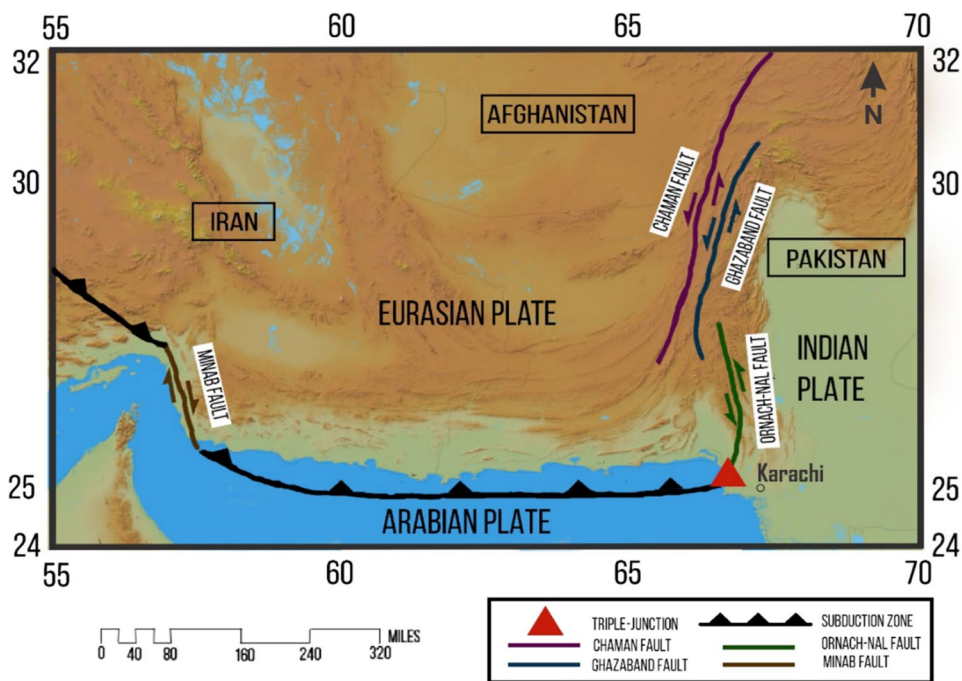
✉ Muhammad Jahangir Khan  
mjahangir.bukc@bahria.edu.pk

<sup>1</sup> Department of Earth and Environmental Sciences, Bahria University Karachi Campus, Karachi, Pakistan

<sup>2</sup> Department of Earth and Environmental Sciences, Bahria University Islamabad Campus, Islamabad, Pakistan

<sup>3</sup> CAS Key Laboratory of Ocean and Marginal Sea Geology, South China Sea Institute of Oceanology, Guangzhou, China

**Fig. 1** Map of Makran onshore (Pakistan and Iran) region. Satellite image by ALOS Global Digital Surface Model (Japan Aerospace Exploration Agency) is modified by adding tectonic features (fault lines and subduction zone) following Regard et al. (2010)



Although it is very rigorous effort but crucial to study the seismological behavior of individual faults particularly identifying their active segments, locking depth of each strike-slip/transform fault within complex Makran region, the seismotectonic behavior can help to delineate the nature of fault blocks either creeping or locked in active Makran subduction region (Szeliga 2012) which may correspond to evaluate the probability of future seismicity hazard in the Makran/coastal belt of Pakistan (Ali and Khan 2015). These boundary faults are accommodating seismic stresses being generated in active Makran subduction system (MSS) (Fattahi and Amelung 2016; Hadi et al. 2013).

The Chaman fault system runs about 900 km north–south along western margin of Indian Plate in MKOR, attributes to seismogeological characteristics and structural heterogeneities in Pakistan, Iran and Afghanistan (Ambraseys and Bilham 2014; Bilham et al. 2007). The GF runs parallel to the CF, extends about 300 km along the north–south direction (Hadi et al. 2013; Smith et al. 2013) and connects with the thrust faults of Sulaiman range near Quetta, Pakistan (Yeats et al. 1979). ONF lies to the south of CF and runs 160 km along the north–south direction toward the Arabian Sea and seems to connect with the Murray Ridge–offshore Pakistan. (The MF which marks western boundary of Makran accretionary prism in the south runs on the Arabian and Eurasian plates boundary about 50 km parallel to the Strait of Hormoz, Iran (Barnhart et al. 2014; Regard et al. 2005.) The tectonic earthquakes in MKOR (Pakistan and Iran) may have primitive feeding from the sinking Arabian Plate at Makran subduction zone and grinding margins of the Indian and Eurasian plates at western margin of Indian Plate and

the transform/transpressional faults in the MKOR which are more susceptible to release the seismic stresses (Regard et al. 2010; Ambraseys and Bilham 2003).

Makran has a long history of earthquakes including the oldest tsunami generated after an earthquake reported 325 B.C (Quittmeyer and Jacob 1979). There were several earthquakes striking the Karachi coast accompanied by tsunamis, such as 1914 northern Makran earthquake, 1945 Pasni earthquake, 1984 offshore Makran/Murray ridge earthquakes and 2002 Ormara earthquake (Khan 2015). The study area has experienced some of the most disastrous earthquakes, e.g.,  $M_w$  7.7 Quetta (1935),  $M_w$  8.1 Makran (1945) and the most recent  $M_w$  7.7 Awaran (2013) earthquake. Among these major events, the Quetta earthquake of 1935 and Awaran earthquake of 2013 occurred along the Ghazaband and splay of Chaman fault (i.e., Hoshab fault). The translational motion between the Indian and Eurasia plates on an active left-lateral Chaman fault system has historically generated some major earthquakes destroying Kabul, Afghanistan and surrounding areas in 1505 (Barnhart et al. 2014). The Quetta earthquake is the deadliest earthquake of southern Pakistan, killed about 35,000–60,000 people (Ambraseys and Bilham 2003). Makran earthquake of  $M_w$  8.1 in 1945 (Byrne et al. 1992) is suggesting the potential of subduction zone to generate earthquakes along the megathrust.  $M_w$  7.7 Awaran earthquake of 2013 is the largest magnitude earthquake of recent decade (Khan 2015). Smith et al. (2013) inferred that the shallow dipping long thrusting Arabian Plate is under a thermal maturation process which may undergone regional extended rupture and can trigger an earthquake of  $M_w$  8.7–9.2. Khan (2015) did integration of focal mechanism

solutions with 3D seismic data (Makran offshore) and studied the involvement of thrust faults in offshore Makran and deep structural seismic activity of MSS. The US geological survey catalog provides dataset of more than a thousand tectonic earthquakes which hit the region over a century (1902–2015) (Ali and Khan 2015), though an updated comparative seismicity response of the individual fault needs to be studied (Crupa et al. 2017; Zinke et al. 2014).

## Data and method

Primary dataset employed in this study includes (1) geographically distributed earthquake data, (2) geological parameters of causative faults, (3) satellite imageries (ASTER GDEM) and (4) published structural maps of MKOR. The point data of earthquakes variables (date, moment magnitude, focal depth, time and epicenter location) were used for seismicity plotting, analysis and evaluation of the understudy faults. The geological characteristics (strike, dip, and rake) of causative faults were retrieved to study the colinear earthquakes above  $M_w$  5.0. The surface geological faults of MKOR (CF, ONF, GF, MF) were traced and georeferenced for comparative seismicity evaluation from structural maps (Ali and Khan 2015; Barnhart et al. 2014; Regard et al. 2010).

We have queried the national and international/open source agencies, e.g., USGS/IRIS—National Earthquake Information Center, Harvard CMT- Lamont-Doherty Earth Observatory (LDEO), for necessary datasets from their archives, since the earthquake recording has started in 1970s in MKOR thus enabled us to consider the instrumental data after 1978. Although the catalog is being updated since Nov 15, 2019, to be latest catalog of the MKOR, a catalog of earthquake was prepared, containing 1015 events from January 1, 1978, to December 12, 2018 (40 years), spatially distributed over MKOR (20°–32° N and 55°–70° E). The catalog was unified to  $M_w$  by adopting empirical relations to convert other magnitudes, i.e., mb, ms &  $M$  to  $M_w$  as derived by Khan et al. (2018) and Scordilis (2006). We have explored the Zmap application developed by Wiemer (2001) to utilize for earthquake data filtering/processing, i.e., declustering of the independent events, historical analysis and magnitude of completeness. The catalog was declustered by following the standard Reasenbergl Declustering method. The declustering operation found 23 clusters of after and foreshocks earthquakes, and a total of 172 events (out of 1033) were removed from the initial data. This exercise presented the catalog of 884 events for further analysis and study. The events of declustered catalog range in magnitude (moment magnitude) from 4.1 to 7.7 and depth 2 to 183 km. Most recent events of the catalog were an intraplate event of  $M_w$  4.4 occurred on August 23, 2019. Zmap tools such

as spatial data function  $F(x,y)$  tools were utilized to estimate the magnitude of completeness ( $M_c$ ). There are different procedures to estimate  $M_c$  such as at max curvature, at fixed  $M$ , at Best 90, 95, etc. We have considered the max. curvature solution to determine the  $M_c$  of catalog. However, it is observed that  $M_c$  highly variable based on the statistical procedure used to estimate it, but the results remain unchanged when choose “Mc 90” or “McBestCombo.”  $M_c$  describes the “the lowest magnitude at which 100% of the earthquakes in a space-time volume are detected” which is useful to reduce the uncertainty about the completeness of the catalog (Shi and Bolt 1982). The  $M_c$  estimations play integral role in  $b$ -value determination. Moreover, the return period of a fault can be well-projected with correct estimates of  $M_c$  and corresponding  $b$ -value.

Geospatial information system (GIS) and satellite remote sensing images are facilitated in integrated seismicity analysis and modeling of spatially distributed earthquakes associated with regional faults. The satellite images of ASTER GDEM were used to trace the fault lines and overlaid seismicity parameters, to understand the regional surface topography along the fault lines and to provide 3D view of focal depth distribution. The ASTER GDEM provided a base map for seismicity analysis. The surface area of CF and GF was covered by 24 tiles (each tile covers  $1^\circ \times 1^\circ$ , horizontal resolution 75 m), while ONF and MF were covered by six tiles. The GDEM tiles of under-investigation regions were merged to build a mosaic (which served as a platform to signify earthquake parameters in relation with the geological fault lines) and to understand the regional surface topography along the fault lines. The mapping applications (ArcMap and ArcScene) were used for seismicity parameters analysis (spatial data mapping, iso-parameters interpolation), preparing the spherical models of fault mechanism solution (FMS), georeferencing of faults lines and classifying the earthquake source parameters (focal depth, magnitude, epicenter locations) to be presented in map layouts.

## Results

The results of this study are summarized into (1) seismological profile of boundary faults and (2) the estimation of re-occurrence period of the potential faults.

### Seismological profile of faults onshore

#### Seismicity of Chaman fault (CF)

The sinistral Chaman fault system (CFS), which represents the western margin of Indian Plate in MKOR, is the largest fault system of Pakistan (Fattahi and Amelung 2016)

and forms a transpressional boundary between the Indian and Eurasian Plates (Hadi et al. 2013). The CF is terminated into Herat fault to the north (near Pak-Afghanistan border) and is branched into multiple curved faults to the south in MKOR. The most damaging earthquake ( $M_w$  6.5) occurred on CF in 1892 near the Chaman town. After this earthquake, the fault was given name CF (Griesbach 1893). Various valuable studies including geological mapping of sediments deformation/slippage, GPS and InSAR investigations were conducted to find out the quantum of elastic stresses being accumulated within the CFS system (Hadi et al. 2013). Recent studies have suggested 8.5 mm/yr creeping rate along CF, which accounts for approximately 30% of the slip rate between the Eurasian and Indian plates (Szeliga et al. 2012). These studies have also pointed out some estimates of slip rates: 18.1 mm/year through sporadic GPS (Mohadjer et al. 2010) and slow slip rate of approximately 8 mm/yr through InSAR analysis (Barnhart et al. 2014). The slip rate of CFS is 19–24 mm/year over 20–25 My based on Khojak Flysch deposit offset [Eocene–Oligocene–Miocene] (Lawrence et al. 1992). Seismicity of Afghanistan and northwestern Pakistan is usually associated with active fault system related to the CF through Herat fault (Furuya and Satyabala 2008; Yeats et al. 1979). During the last four decades (1978–2018), seismicity along the CF has low magnitude  $M_w$  4–5 (Table 1). The largest earthquake occurred at CF was  $M_w$  5.7 in 1978. The earthquakes magnitude classification in various parts of CF is shown through different colors and radii of circles in Fig. 2a. There are 23 colinear epicenters which overlaid the CF (18 jolts higher than  $M_w$  4, and 5 earthquakes higher than  $M_w$  5). The focal depths of these events are used to create the subsurface planner view of focal depths created through interpolation (Fig. 2b). The iso-focal depth map of earthquakes lies in the buffer of CF and reveals that the middle section is related to multiple peaks of focal depth which may be interpreted as shallow subsections of the deformed CF. CF surface view is forming a two-way gentle dipping plans in the subsurface under the surface fault trace of CF, i.e., the focal depths increase toward north and south. The FMS of  $M_w > 5$  events unveils the fault nature and slippage of blocks along the CF. The FMS of CF consists four quadrants and exhibits strike-slip

movement of the blocks (Fig. 3). It is apparent that the dip of the fault plan is toward north and responsible for more transpressional stress (Fig. 3). The recurrence period of CF is estimated in following section for intermediate magnitude events, i.e.,  $M_w$  4–5, which suggests that CF may trigger another event in 9 years (Table 2). CF is most vulnerable seismotectonic element of Pakistan. Surface fractures and ground displacement are evident from historic events along CF (1505 Kabul earthquake, Babur 1912). Griesbach (1893) documented that 1892 earthquake in Pakistan caused bending of iron made rail track by 0.75 m, and  $M_w$  6.1 Naushki earthquake ruptured the surface.

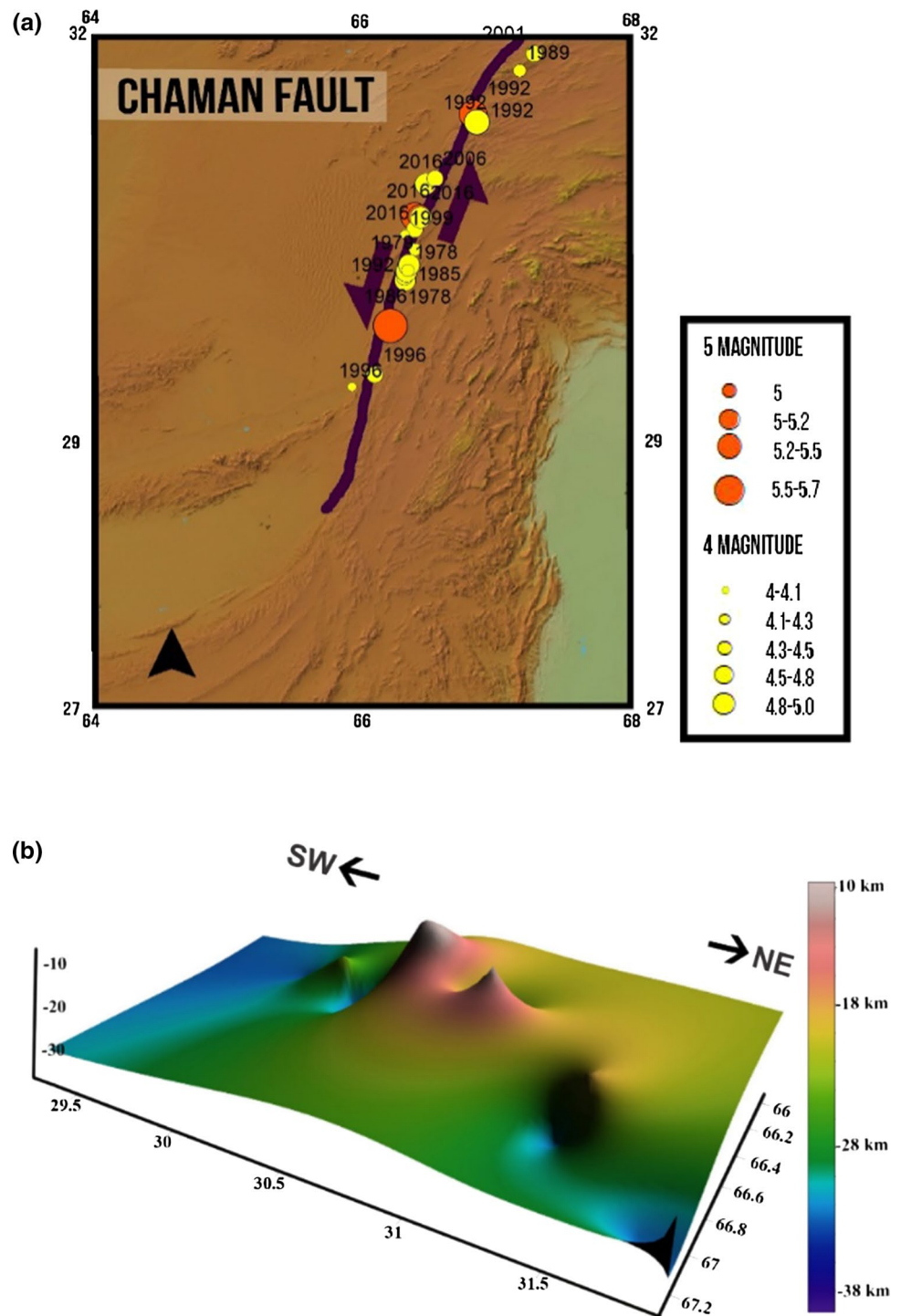
### Seismicity of Ghazaband fault

GF accommodates strike-slip stress deformation in the MKOR between 27° and 31° N (~40 km east of Chaman town, Balochistan) and has a greater potential of seismicity. Slip along the GF plane was responsible for the extensive damage to Quetta in the 1935  $M_w$  7.7 earthquake (Byrne et al. 1992; Ramanathan and Mukherji 1938). This was the deadliest earthquake of the region, killing about 35,000–60,000 people and injuring several thousands. The GF was recognized after this earthquake as the most seismically active part of CFS. Recent research (Szeliga 2012) reported reliable estimate of combined velocity at CF and GF, determined through GPS network (30°–32° N), is approx. 12 mm/year, while the Indian Plate is moving toward Eurasia (Afghan block) with a velocity of 24–28 mm/year (Altamimi et al. 2011). In the CFS block, velocity increases toward GF from 6 to 8.5 mm/year (Szeliga 2012). The GF faced a slip of approx. 9 cm from left-laterally displaced shallow fault from a  $M_w$  5.5 earthquake occurred on October 2007 (Fattahi et al. 2015). During the last four decades (1978–2018), 38 earthquakes of varied magnitude have occurred on GF which are given in Table 1. The GF has ranked as the highest seismicity producing fault among the understudy faults of MKOR. Figure 4a demonstrates the earthquakes of different magnitudes lying on the GF surface traced over GDEM. The focal depths of these events are used to create the subsurface planner view of focal depths through the interpolation (Fig. 4b). The focal depths of the events lie within range of 10–33 km, which means the deformation of the GF is relatively shallow and fast as compared to adjacent CF. GF generated relatively more events within 0–10 km. The iso-focal depth plane of GF is forming an asymmetrical pattern of alternate crest and trough suggesting highly deformed blocks (slip along z -axis) with some horizontal displacement ascertained in strike-slip deformation along the GF plane. FMS of earthquake higher than  $M_w$  5.0 along the GF divulge the fault nature and slippage along the GF (Fig. 3). The spheres of FMS are four quadrants exhibiting the shear stresses of the blocks.

**Table 1** Number of earthquakes occurred along the understudy faults

Magnitude range	CF	GF	ONF	MF
4.0–4.4	12	16	10	9
4.5–4.9	6	13	12	7
5.0–5.4	2	6	1	–
5.5–5.9	3	2	3	1
6.0–6.4	0	1	0	0
6.5–6.9	0	0	0	0

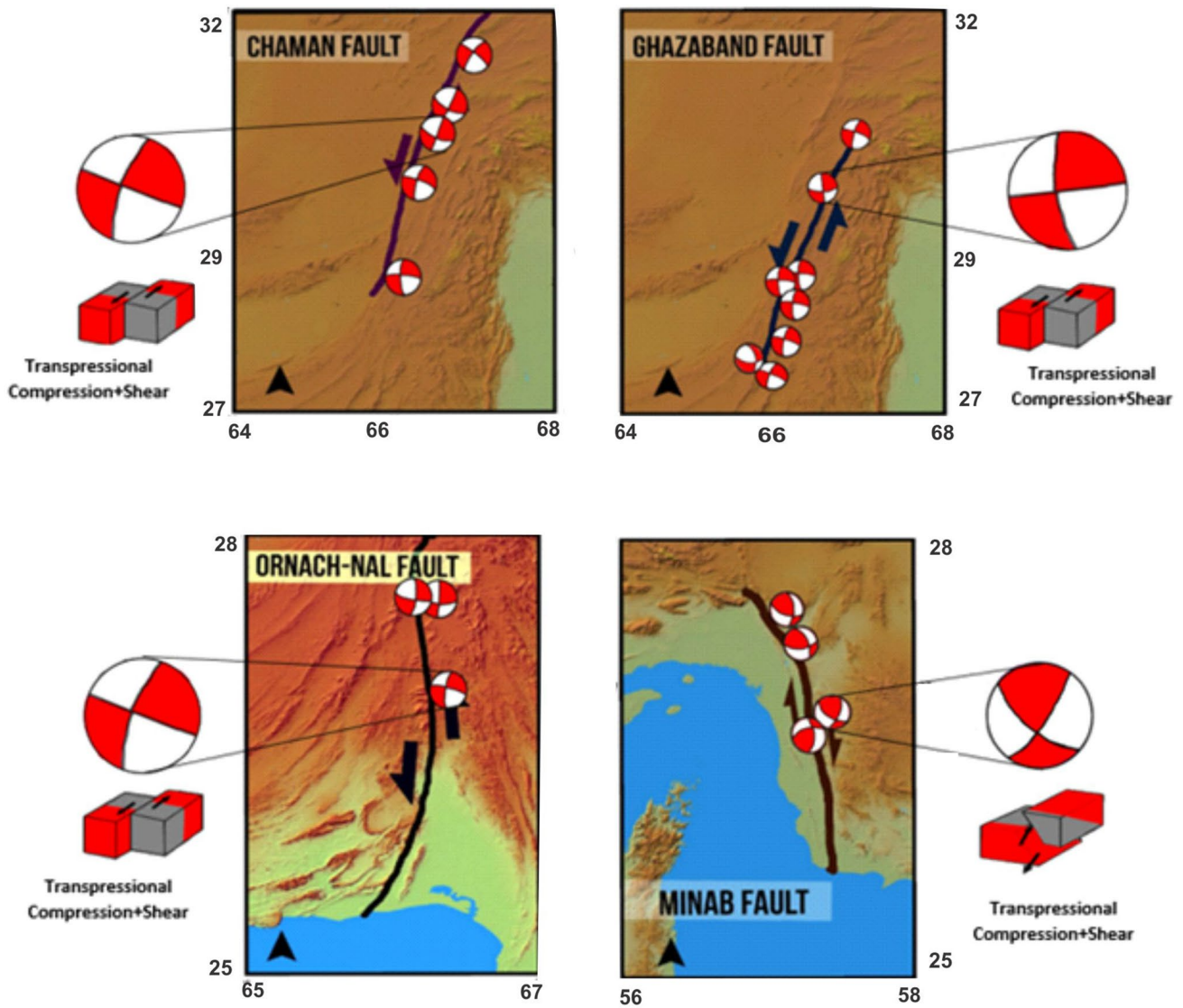
**Fig. 2** **a** Earthquake data displayed over the georeferenced Chaman Fault, exported from ArcMap. The DEM surface depicts the surface topography along the CF. Legend displays the earthquake magnitude through different color and size. **b** Iso-focal depth of earthquakes occurred along the Chaman Fault, exported from ArcMap. Legend displays the focal depth in km classified in various color bars



### Seismicity of Ornach-Nal fault

The ONF extends in MKOR from 25.5 to 28°N. The sinistral motion elongates about 250 km in southernmost onshore region along the ONF (Szeliga 2012). The ONF is associated with a triple junction of Arabian, Eurasian and Indian plates which are located at the south of Somiani Bay and connects the Murray Ridge/Makran subduction zone. The literature

review suggests that the slip rates of ONF are roughly 20–40 mm/year (Lawrence et al. 1992) The average velocity measured on the triple junction south of ONF is 15.1 mm/year (Szeliga 2012). The seismicity along the ONF is always noteworthy being closer to metropolitan city Karachi (east of ONF and the triple junction) and connecting onshore and offshore. It is observed that the ONF region was seismically quite before 1972 (International Seismological Center



**Fig. 3** Focal mechanism solutions (Red and white shaded spheres) of events  $M_w > 5$  magnitude displayed on Chaman fault, Ghazaband fault, Ornach-Nal fault and Minab fault (Source of focal mechanism

parameters of earthquakes occurred along understudy faults was Global Centroid Moment Tensor Project Dziewonski et al. 1981; Ekström et al. 2012)

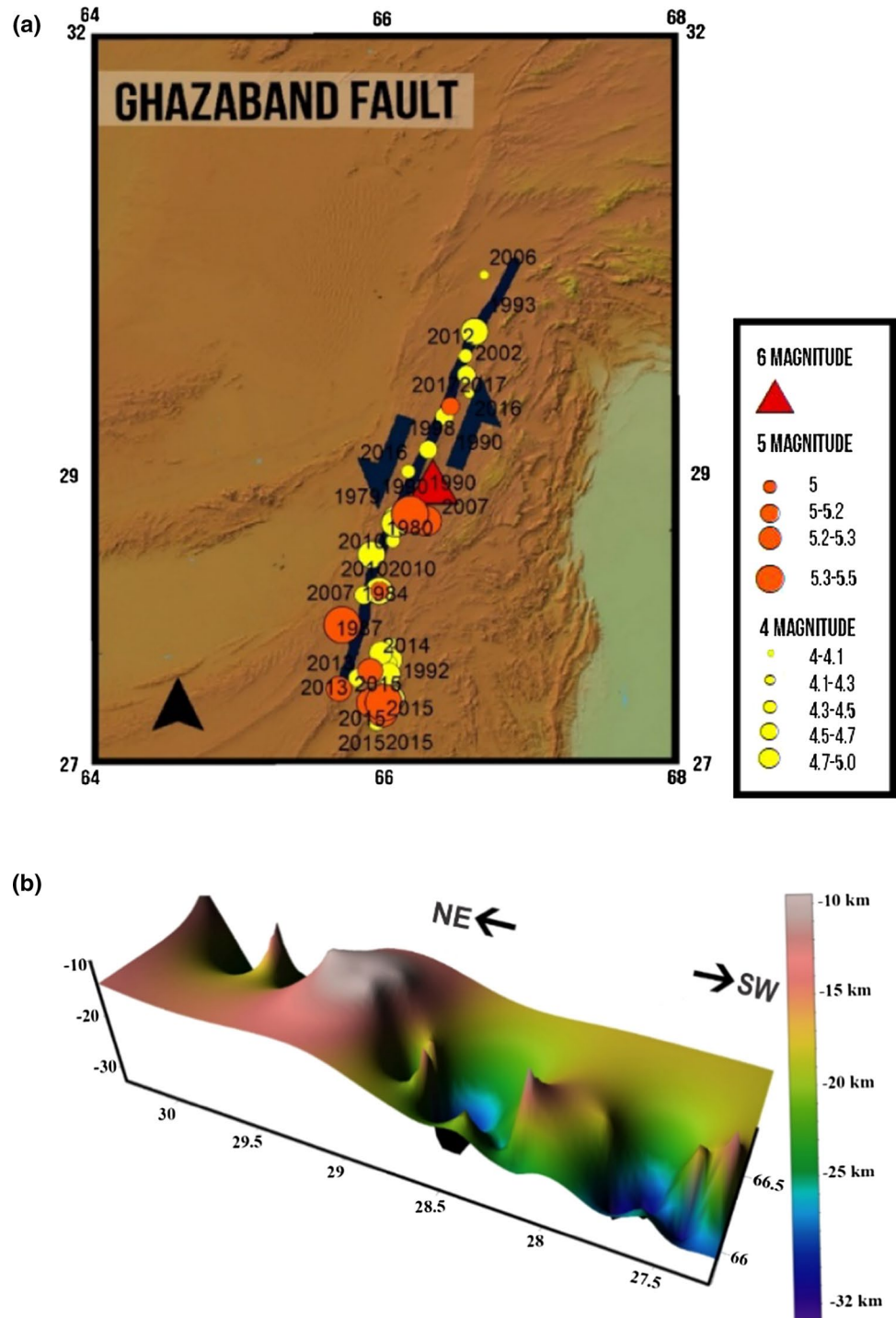
**Table 2** Estimations of recurrence interval for moderate events at corresponding fault

Faults	$T^{\circ}$ (years)	$b$ -value	$M^*$	$M$	$N$	Tr (year)
CF	40	$0.457 \pm 0.191$	4	5.7	23	$9 \pm 1$
GF	40	$0.644 \pm 0.074$	4	6.1	38	$14 \pm 1$
ONF	40	$0.530 \pm 0.404$	4	5.5	26	$07 \pm 0.5$
MF	40	$0.666 \pm 0.155$	4	5.5	17	$15 \pm 0.8$

Dataset) either it is a delusion or non-recording era. After 1970s, the seismicity observed at ONF nevertheless of small magnitude. This study highlights 23 earthquakes with  $M_w$  4–4.9, 3 earthquakes with  $M_w$  5–5.9 and the highest  $M_w$  5.6, occurred on ONF during the last four decades (Table 1). The earthquake magnitude distribution is classified (Fig. 5a) and

overlaid the georeferenced ONF on GDEM Image. The ONF is a crust cutting fault as it has generated more deeper focus earthquakes up to 42 km (Fig. 5b). It is envisaged that there could be the involvement of the subducting Indian Plate section or the underthrusting of the Arabian Plate to this depth. The iso-focal depth visualization follows a similar character

**Fig. 4** **a** Earthquake data displayed over the georeferenced Ghazaband fault over DEM surface. Legend displays the earthquake magnitude through different color and size. **b** Iso-focal depth of earthquakes along the Ghazaband fault, exported from ArcMap. Legend displays the focal depth in kilometers classified in various color bars

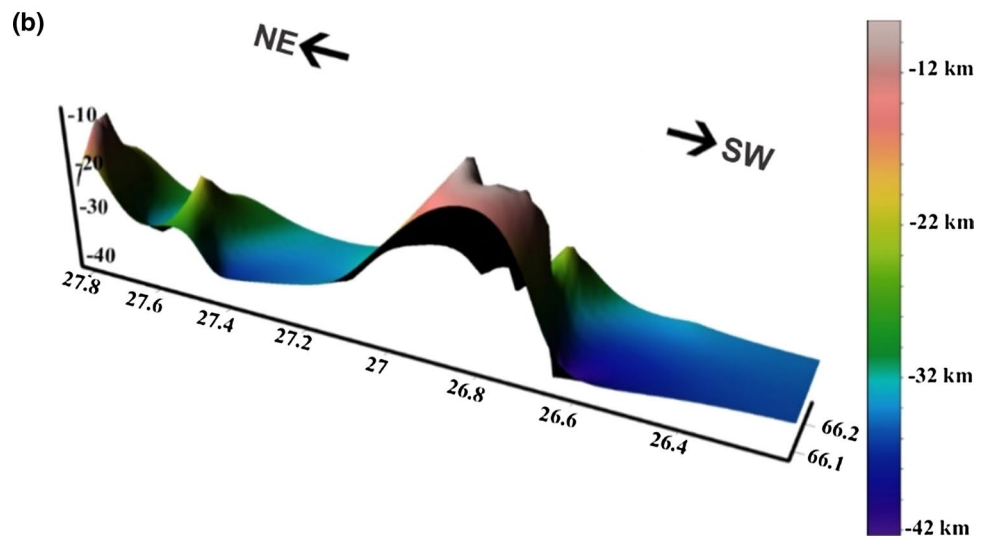
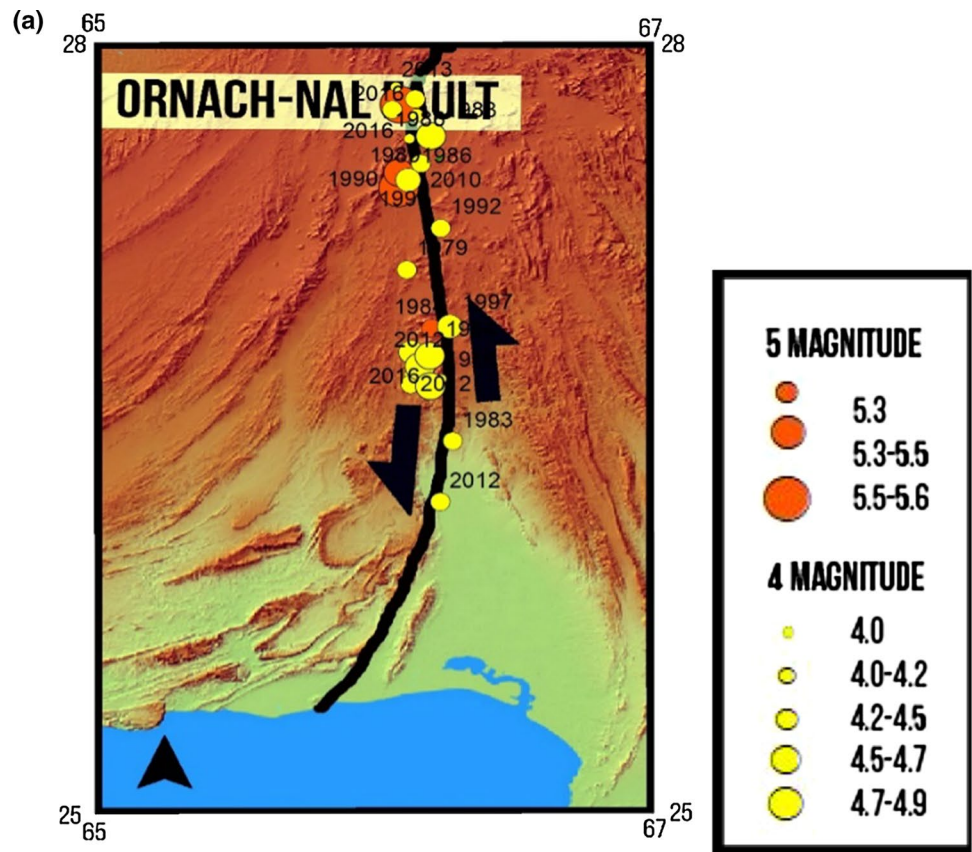


of the CF, i.e., two-way dipping, northward and southward; however, the middle portion of ONF has generated earthquakes from shallow depths. The undulating depth profile (Fig. 5b) of ONF plane highlights the horizontal and vertical displacements of the ONF sections. The FMS of  $M_w > 5$  reveals the nature of block movements of ONF (Fig. 3) and shows transpressional stress. The earthquakes along ONF have deep focal depth (i.e., crust cutting events) than the

other faults. Does Murray Ridge connect with ONF? It is suggested that a detailed study of intermediate to deep focal depths and seismic reflection images (synthetic seismic survey for E&P sector) may shed some light on the connection. The mechanism of left-lateral motion (onshore) and right lateral motion appeared along the Sonne fault (Kukowski et al. 2001) and Owen fracture zones (offshore Arabian sea) which also needs to be resolved (Yeats et al. 1979).



**Fig. 5 a** Earthquake data displayed over the georeferenced ONF over GDEM surface. Legend displays the earthquake magnitude through different color and size. **b** Iso-focal depth of earthquakes occurred along the Ornach-Nal fault, exported from ArcMap. Legend displays the focal depth in km classified in various color bars



The estimate of recurrence interval on ONF has approx.  $707 \pm 0.5$  years for  $M_w$  5.5.

**Seismicity of Minab fault**

MF is bounding the MKOR from west and represents the seismicity of western Makran onshore. MF connects the Zagros mountains and Makran subduction zone. The

western Makran neighbors (strait of Hormoz and Zendan-Minab fault) are seismically active (Regard et al. 2005); however, global seismic networks and historical records indicate that the seismic activity is low, and earthquakes depth is also shallow across the area around MF (Nemati 2019). MF is located near the coast, perpendicular to the shoreline, and seems to have penetrated into Makran subduction zone. It extends for about 50 km and runs parallel

to the Strait of Hormoz (Regard et al. 2010). MF is 25°–45° dipping discontinuous and nonlinear thrust faults affecting the Late Cenozoic and the Quaternary geological deposits closely associated with the Minab fold (Burg 2018). The Arabian Plate is moving to the direction of Strait of Hormuz with a velocity of 23 mm/year (Ahmed et al. 2018). The Zagros thrust belt is moving in southeastern end up to 10 mm/year. The average strike-slip rate of Minab–Zendan fault is 12 mm/year (Peyret et al. 2009). The seismicity source parameters associated with the MF (1978–2018) are overlaid with fault trace of MF. The spatial distribution of the magnitude highlights the 16 earthquakes with  $M_w$  4–4.9 and 1 earthquake with  $M_w$  5.5 occurred on MF (Table 1, Fig. 6a). The MF has lowest seismicity potential among all the understudy faults. The iso-focal depth (Fig. 6b) exposed the shallow deformed middle portion of the MF plane with a sagging block/basin (in the northern MF). The associated FMS of  $M_w > 5$  along the MF reveals dominant involvement of the reverse faulting and shear strain. The oblique fold cross-cuts MF in the south while the thrust splits two well-defined structures in the northern MF. The western part is an east dipping reverse fault affecting the late Pleistocene front of the fold and thrust belt fans, while the eastern part is an inactive northeast oblique dipping thrust that cross-cuts the MF (Nemati 2015; Regard et al. 2005). The estimated recurrence period of MF is  $1515 \pm 0.8$  years (Table 2).

## Recurrence interval estimation

The historical data analysis of the catalog was performed, in addition to collection of historical events of the region from the exhaustive catalogs of Ambraseys and Bilham (2002), and published studies of Smith et al. (2013), Ali and Khan (2015), Khan et al. (2018). Since the seismological network was established in 1970s, the information of the events pre-1970s uncertain the precise magnitude, focal depth and location in the complex geological environment. It is apparently found that those historic events with strong impacts/damageability were considered probably due to their tangible impacts on property and life. However, lack of seismological knowledge and scientific background of the historian or people settled in remote tribble belts of Balochistan (Pakistan) and Sistan (Iran) preclude exhaustive information of historic events in MKOR. Still, the region is lacking in richness of seismological stations and state-of-the-art instrumentation to monitor static stress in the faulty blocks of Makran onshore region (MKOR) and early warning tsunami sensors along the Makran coast, etc. (Ellouz-Zimmermann et al. 2007; Kukowski et al. 2001; Minshull et al. 2015; Frohling and Szeliga 2016).

The understudy faults are originating earthquakes along their lineaments. The estimation of the recurrence of

seismicity along CF, GF, ONF and MF is done by Eq. 1 of Ali and Khan (2015).

$$\text{Tr} = (T^\circ \times 10^{b(M-M^*)})/N \quad (1)$$

Tr = Recurrence Time Period in years,  $T^\circ$  = Observational Time Period in years,  $M^*$  = threshold magnitude,  $N$  = cumulative number of earthquakes of magnitude  $M^*$  and above.  $M$  = maximum credible magnitude.

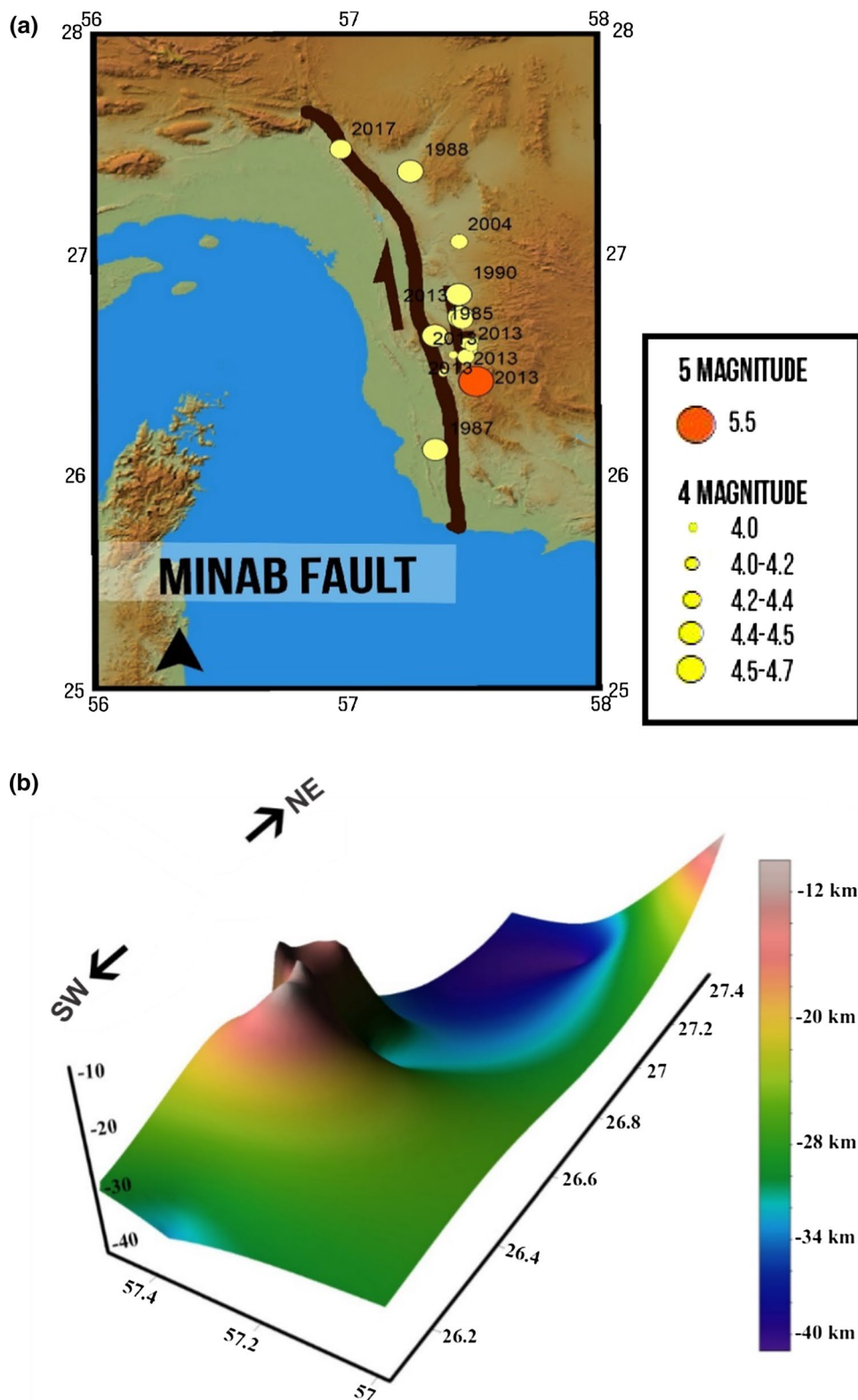
The spatial distribution of  $b$ -values in southern Pakistan is shown in Fig. 7, which help to deduce the distribution of  $b$ -values in MKOR and key faults. Although some temporal and spatial variations in  $b$ -values are also reported in MKOR region (Rani et al. 2011; Ali and Khan 2015), the estimates of  $b$ -value for specific regions of CF, GF, ONF, and MF are abstracted in Table 2. There is a significant low  $b$ -value for CFS segment which generated events from 1978 to 2019. This is an anomalous situation in the light of the previous literature which attributes low- $b$  values with high seismic probability. Thus, lowness of  $b$ -values in CFS segment might be treated as a precursor of future event (Shi and Bolt 1982; Smith 1981).

The researchers (Crupa et al. 2017; Fattahi and Amelung 2016; Barnhart et al. 2014; Szeliga et al. 2012; Regard et al. 2010) utilized satellite images, e.g., interferometric synthetic aperture radar (InSAR) to infer quantitative information about the ground deformation rates, along with locking depths of the CF and GF segments. The GPS- and InSAR-based plate kinematics reported that the plates movement measured near 26° N across the ONF ranges between 13.4 and 16.9 mm/year (avg. 15.1 mm/year) within less than 3 km locking depth (Fattahi et al. 2015; Szeliga 2012). The velocity measured at the town of Chaman (near 30° N) is 8.5 mm/year (6.8–10.3 mm/year), and the CF is locked at ~3.4 km depth (Szeliga 2012). The creeping block velocity varies at places, e.g., across CF it is 14.1–19.5 mm/year, and the convergence rate is approx. 2 mm/year near western most strike-slip fault and 6–9 mm/year near the transpressional faulting (Frohling and Szeliga 2016; Altamimi et al. 2011). The different tectonic stresses associated with Indian, Arabian and Afghan blocks are being accumulated at tectonic margins due to different geodynamics and megathrust block velocities in the study area (Khan et al. 2008; Hussain et al. 2002).

## Discussion

MKOR has been undergone diversified plate tectonic processes such as the subduction, underthrusting, transpressional faulting and mountain building in geological times. In the west of Pakistan and southeast of Afghanistan, the oblique motion of Indian Plate relative to the Eurasian Plate has resulted in a

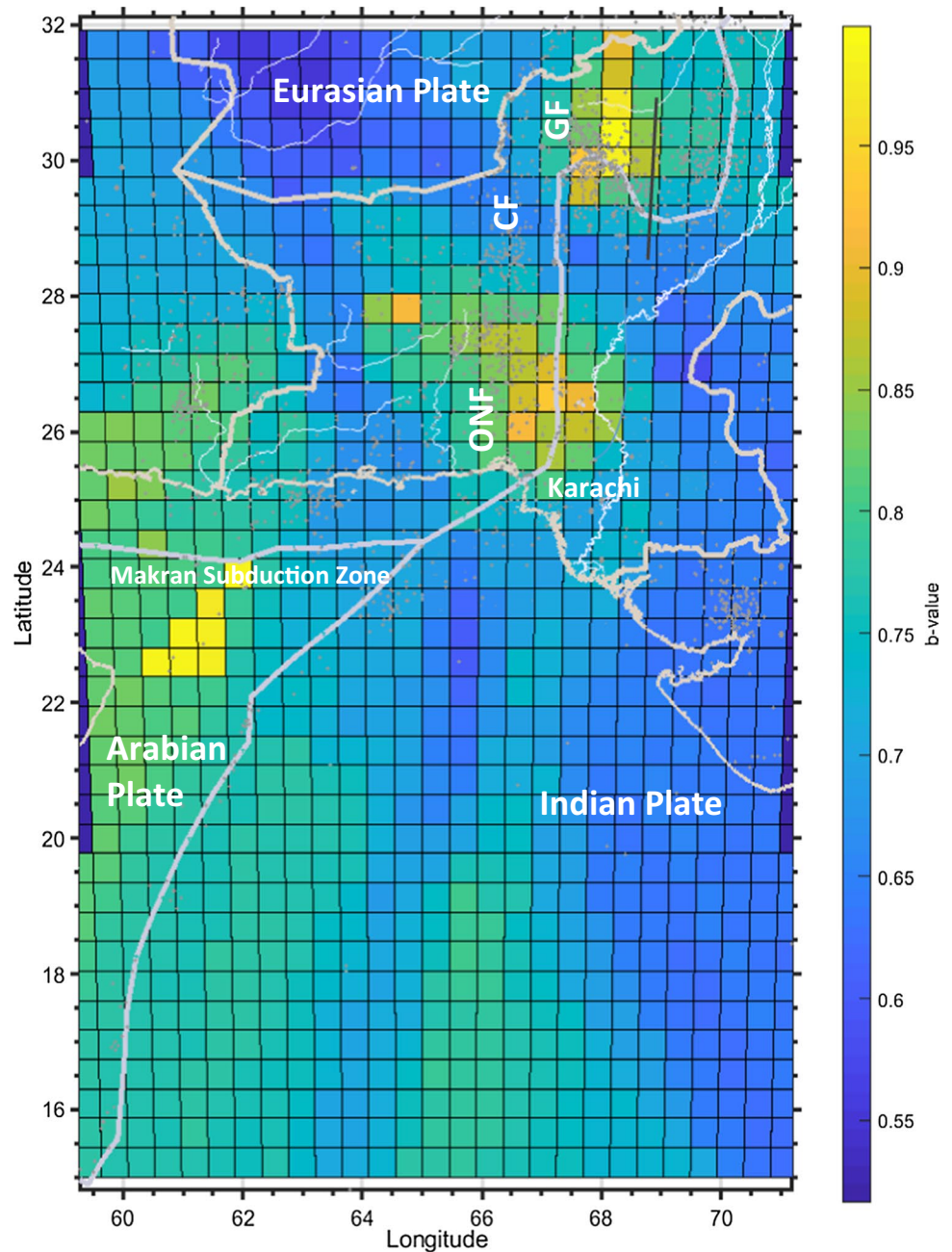
**Fig. 6** **a** Earthquake data displayed over the georeferenced MF over GDEM surface. Legend displays the earthquake magnitude through different color and size. **b** Iso-focal depth of earthquakes occurred along the Ornach-Nal fault, exported from ArcMap. Legend displays the focal depth in km classified in various color bars



complex fold and thrust belt of Sulaiman range (Kazmi and Jan 1997). Makran (onshore) region shows an active seismicity which is linked with the interaction among three major plates (Arabian, Eurasia and Indian). The tectonic outcome in MKOR

exposes on surface the prominent imbricate thrust slices, segmented mountain topography and a dissimilarity in seismogenic behavior in western and eastern Makran as evident by the respective key faults (CF, GF, ONF and MF). Frohling and

**Fig. 7** Spatial distribution of  $b$ -value in southern Pakistan. The  $b$ -values are computed at a grid of 1053 pixels considering  $N=250$  events and max. 50 events less than the Mc. The color ramp indicates the contrasting variation in  $b$ -value in the study region. Smaller regions did not contain enough data for adequate resolution of the  $b$ -values computation raster particularly at margins (left and right sides—deep blue color)



Szeliga (2016) and Penney et al. (2017) investigated elastic strain accumulation in western Makran and expressed that historical records are not long or reliable enough to expect a great event in future. Smith et al. (2013) investigated thermal modeling of the subduction zones and concluded that a potential  $M_w > 9$  earthquake is possible if full length of MSZ ruptures. The literature review (Ali and Khan 2015; Ambraseys and Bilham 2003a, b), past earthquake reports and residents' observations suggest that the regional intraplate earthquakes either occurred in onshore Makran or offshore Makran region [such as Dalbandin (400 km from Karachi), Pasni (500 km away), Awaran (210 km away from Karachi)] have proportionally

jolted Karachi. Any significant earthquake in Makran and its contiguous coastline may have awful impact on populous city Karachi (Smith et al. 2013; Sarwar and Alizai 2013), thus, a significant regional hazard mitigation plan in coastal Pakistan (Barnhart et al. 2013; Smith et al. 2013; Hatzfeld and Molnar 2010). October 8, 2005, Kashmir earthquake gave a spark to seismological studies of Pakistan which categorized the region adjacent to CF and GF as active seismotectonic zones, whereas ONF lies in the least seismicity zone (Waseem et al. 2019). However, recent studies (Barnhart et al. 2013; Smith et al. 2013; Szeliga et al. 2012) and the existence of nearby triple junction of Eurasian, Arabian and Indian plates suggested

that entire Makran shall fall within high-risk seismotectonic zone of Pakistan. Although the earthquakes  $M_w$  7.7 Dalbandin (January 18, 2012) and  $M_w$  7.7 Awaran (September 24, 2013) have released most of the elastic strain energy accumulated in MSS, the remaining energy in the system under the influence of active tectonics of triple junction may trigger a major tremor affecting the Balochistan and Sindh. Hence, the recommendation of redesigning the building codes for proximal coastal cities of Karachi and Gawadar would be a safety measure against any disaster (Ali and Khan 2014).

Considering the short history of earthquake observations compared to the potential length of a seismic cycle, we cannot exactly foresee the next rupture, when and where it can occur in MKOR. The uncertainty in evaluated seismic hazard is due to moderate and large plate margins grinding the crust in various directions. Moderate thrust, strike-slip or transpressional motion can generate earthquakes in shallow depth < 10 km, whereas moderate and large earthquakes, in intermediate depth range (> 30 km), might occur along the under-investigation fault lines of MKOR. The comparative seismicity evaluation in this study contributes to understand the seismicity response of the major faults. It is apparent that GF has higher seismicity than the other faults because it accumulates more elastic stress, thus poses more risk than the other faults. The cities located in vicinity of CF and ONF closely associated with GF are also at risk. The statistical analysis showed that the MF generated most earthquakes in 2013, including an earthquake with  $M_w$  5.5. The Hoshab fault recently generated a  $M_w$  7.7 earthquake, which affected the seismicity of all the faults in the region. In the region of CF,  $M_w > 7$  earthquakes have not been reported previously. Similarly, the ONF have not produced earthquake with  $M_w > 5.9$ . The GF has been triggered by two major events (1931  $M_w$  7.2 Mach earthquake, and 1935  $M_w$  7.5 Quetta earthquake). The focal depth of earthquakes at ONF is roughly 42 km which is deeper than other faults. The ONF is observed as another fault system with different geometry of fault trace, FMS attributes, focal depth profile and swinging splays to the western side of the main fault traces. ONF is closer to the coast and offshore region, with potential of generating  $M_w > 5.6$  may cause a future threat of earthquake vibrations reaching to Karachi. This study highlights the insight of ONF seismicity potential, which may be considered for building codes of Karachi. The FMS provides an insight of the strike-slip nature of all the faults, with a slight influence of compression and subducting Indian Plate.

## Conclusion

This study provides a seismicity comparison between key strike-slip faults of Makran, i.e., CF, GF, ONF and MF over last four decades which can help to evaluate the potential

faults response. The GF is ranked as frequent and exhibiting high seismicity response than other faults. The neighborhood region of GF elucidated by high  $b$ -values in southern Pakistan. The middle section of CF generated shallow earthquakes than those of northern and southern segments. The FMS of CF, GF, ONF and MF consist four quadrants and exhibit dominant component of transpressional nature, certainly under the influence of collisional and/or subduction regime. MF originated seismicity infrequently; however, a major increase in seismicity was observed in after Awaran earthquake. The earthquakes along ONF have deep focal depth (i.e., crust cutting events) than the other faults. It can be postulated that  $M_w$  7.0 or above earthquake might be another destructive event in MKOR; thus, an extensive probabilistic seismicity hazard assessment has recommended for future study and safety.

**Acknowledgements** We acknowledge the guidance of Roger Bilham (Department of Geological Sciences, University of Colorado, Boulder, USA) and Rodolfo Console (Istituto Nazionale di Geofisica e Vulcanologia (INGV), Italy) in drafting the manuscript content, result and seismicity analysis. We also thankful to Pakistan Petroleum Limited and administration of Bahria University Karachi Campus (BUKC) for providing the research facility at PPL Post-graduate laboratory, Department of Earth & Environmental Sciences, BUKC. The dataset of this integrated study is a subset of the data resources being utilized in PhD research work of first author of this manuscript.

## Compliance with ethical standards

**Conflict of interest** There is no conflict of interest among the authors.

## References

- Ahmed S, Hassan S, Mehnood K, Maqsood T (2018) Role of Chaman transform boundary fault in the deformation of Eastern Kharan fore-arc basin. *J Himal Earth Sci* 51:75–98
- Ali M, Khan MJ (2014) Invisible fault lines ruin the development of urbanization. In: Proceedings of 1st international conference on infrastructure, management, assessment and rehabilitation techniques (ICIMART' 14). American University of Sharjah
- Ali M, Khan MJ (2015) GIS based study on seismicity of Makran over 100 years. *Int J Econ Environ Geol* 6(2):11–16
- Altamimi Z, Métivier L, Collilieux X (2011) ITRF 2008 plate motion model. *Geophys Res Abstr* 13, EGU 2011–4750
- Ambraseys N, Bilham R (2003a) Earthquakes and associated deformation in North Baluchistan 1892–2001. *Bull Seismol Soc Am* 93(4):1573–1605
- Ambraseys N, Bilham R (2003b) Earthquakes in Afghanistan. *Seismol Res Lett* 74(2):107–123
- Ambraseys N, Bilham R (2014) The tectonic setting of Bamiyan and seismicity in and near Afghanistan for the past twelve centuries. After the destruction of Giant Buddha statues in Bamiyan (Afghanistan) in 2001. Springer, Berlin, pp 101–152
- Babur ZM (1912) *The Babur-Nama in english translated by A S. Beveridge*. Steven Austin, Hertford
- Barnhart WD, Lohman RB, Mellors RJ (2013) Active accommodation of plate convergence in Southern Iran: earthquake

- locations, triggered aseismic slip, and regional strain rates. *J Geophys Res Solid Earth* 118(10):5699–5711. <https://doi.org/10.1002/jgrb.50380>
- Barnhart W, Hayes G, Briggs R, Gold R, Bilham R (2014) Ball-and-socket tectonic rotation during the 2013 Mw 7.7 Balochistan Earthquake. *Earth Planet Sci Lett* 403:210–216
- Bilham R, Lodi S, Hough S, Bukhary S, Khan AM, Rafeeqi SFA (2007) Seismic hazard in Karachi, Pakistan: uncertain past, uncertain future. *Seismolog Res Lett* 78(6):601–613
- Burg KP (2018) Geology of the onshore Makran accretionary wedge: synthesis and tectonic interpretation. *Earth Sci Rev.* <https://doi.org/10.1016/j.earscirev.2018.09.011>
- Byrne DE, Sykes LR, Davis DM (1992) Great thrust earthquakes and aseismic slip along the plate boundary of the Makran Subduction Zone. *J Geophys Res* 97:449
- Crupa WE, Khan SD, Huang J, Khan AS, Kasi A (2017) Active tectonic deformation of the western Indian plate boundary: a case study from the Chaman fault system. *J Asian Earth Sci* 147:452–468
- Dolati A, Burg JP (2013) Preliminary fault analysis and paleostress evolution in the Makran Fold-and-Thrust Belt in Iran. In: Al Hosani K, Roure F, Ellison R, Lokier S (eds) *Lithosphere dynamics and sedimentary basins: the Arabian plate and analogues*. Springer, Heidelberg, pp 261–277
- Dziewonski AM, Chou TA, Woodhouse JH (1981) Determination of earthquake source parameters from waveform data for studies of global and regional seismicity. *J Geophys Res Solid Earth* 86(B4):2825–2852
- Ekström G, Nettles M, Dziewoński AM (2012) The global CMT project 2004–2010: centroid-moment tensors for 13,017 earthquakes. *Phys Earth Planet Inter* 200:1–9
- Fattahi H, Amelung F (2016) InSAR observations of strain accumulation and fault creep along the Chaman Fault System, Pakistan and Afghanistan. *Geophys Res Lett* 43:8399–8406
- Fattahi H, Amelung F, Chaussard E, Wdowinski S (2015) Coseismic and postseismic deformation due to the 2007 M5.5 Ghazaband fault earthquake, Balochistan, Pakistan. AGU Publications, Washington, pp 3305–3312
- Frohling E, Szeliga W (2016) GPS constraints on interplate locking within the Makran subduction zone. *Geophys J Int* 205(1):67–76
- Furuya M, Satyabala S (2008) Slow earthquake in Afghanistan detected by Insar. American Geophysical Union, Washington
- Griesbach C (1893) Notes on the earthquake in Balochistan on the 20th December 1892. *Geol Surv India* 26(2):57–61
- Hadi SU, Khan SD, Owen LA, Khan AS, Hedrick AK, Caffee MW (2013) Slip-rates along the Chaman fault: implication for transient strain accumulation and strain partitioning along the western Indian plate margin. *Tectonophysics* 608(2013):389–400
- Hatzfeld D, Molnar P (2010) Comparisons of the kinematics and deep structures of the Zagros and Himalaya and of the Iranian and Tibetan plateaus and geodynamic implications. *Rev Geophys* 48:RG2005
- Hussain J, Butt KA, Pervaiz K (2002) Makran coast: a potential seismic risk belt. *Geological Bulletin University, Peshawar*, pp 43–56
- Kazmi A, Jan Q (1997) *Geology and tectonics of Pakistan*. Graphic Publishers, Santa Ana
- Khan MJ (2015) *Integrated study on seismicity of Makran over 100 years*. Department of Earth & Environmental Sciences, BUKC. MS Thesis
- Khan MA, Bendick R, Ismail Bhat M, Bilham R, Kakar DM, Faisal Khan S, Lodi SH, Sufyan M, Singh B, Szeliga W, Wahab A (2008) Preliminary geodetic constraints on plate boundary deformation on the western edge of the Indian plate from TriG-Gnet (Tri-University GPS Geodesy Network). *J Himal Earth Sci* 41:71–87
- Khan S, Waseem M, Khan MA, Ahmed W (2018) Updated earthquake catalogue for seismic hazard analysis in Pakistan. *J Seismol* 22(4):841–861
- Kukowski N, Schillhorn T, Huhn K, von Rad U, Husen S, Flueh ER (2001) Morphotectonics and Mechanics of the central Makran accretionary wedge off Pakistan. *Mar Geol* 173:1–19
- McCall GJH (2003) A critique of the analogy between Archaean and Phanerozoic tectonics based on regional mapping of the Mesozoic-Cenozoic plate convergent zone in the Makran, Iran. *Precambrian Res* 127:5–17
- Mohadjer S, Bendick R, Ischuk A, Kuzikov S, Saydullaev U, Lodi S, Zubovich A (2010) Partitioning of India- Eurasia convergence in the Pamir- Hindu Kush from GPS measurements. American Geophysical Union, Washington
- Nemati M (2015) Aftershocks investigation of 2010 Dec. and 2011. Jan Rigan earthquakes in the southern Kerman province, SE Iran. *J Tethys (Iran)* 3(2):96–113
- Nemati M (2019) Seismotectonic and seismicity of Makran, a bimodal subduction zone, SE Iran. *J Asian Earth Sci* 169:139–161
- Penney C, Tavakoli F, Saadat A, Nankali HR, Sedighi M, Khorrami F, Sobouti F, Rafi Z, Copley A, Jackson J, Priestley Keith (2017) Geodynamics and tectonics Megathrust and accretionary wedge properties and behaviour in the Makran subduction zone. *Geophys J Int* 209:1800–1830
- Peyret M, Djamour Y, Hessami K, Regard V, Bellier O, Vernant P, Daignieres M, Nankali H, Van Gorp S, Rigoulay M, Goudarzi M (2009) Present-day strain distribution across the Minab-Zendan-Palami fault system from dense GPS transects. *Geophys J Int* 179:751–762
- Quittmeyer RC, Jacob KH (1979) Historical and modern seismicity of Pakistan, Afghanistan, northwestern India, and southeastern Iran. *Bull Seismol Soc Am* 69, 773–823, p 805, Appendix 3, citing Oldham T (1882) A catalogue of Indian earthquakes from the earliest time to the end of A.D. 1869. *Mem Geol Surv India* 19: 163–215
- Quittmeyer RC, Kafka A (1984) Constraints on plate motions in southern Pakistan and the northern Arabian Sea from the focal mechanisms of small earthquakes. *J Geophys Res: Solid Earth* 89(B4):2444–2458
- Ramanathan K, Mukherji S (1938) A seismological study of the Baluchistan, Quetta, Earthquake of May 31, 1935. Geological Survey of India, Kolkata, pp 483–513
- Rani VS, Srivastava K, Srinagesh D, Dimri VP (2011) Spatial and temporal variations of b-value and fractal analysis for the Makran region. *Mar Geodesy* 34(1):77–82
- Regard V, Bellier O, Thomas J-C, Boul'Es D, Bonnet S, Abbassi M, Feghhi K (2005) Cumulative right-lateral fault slip rate across the Zagros–Makran transfer zone: role of the Minab–Zendan fault system in accommodating Arabia-Eurasia convergence in Southeast Iran. *Geophys J Int* 162:177–203
- Regard V, Hatzfeld D, Molinaro M, Aubourg C, Bayer R, Bellier O, Yamini-Fard F, Peyret M, Abbassi M (2010) The transition between Makran subduction and the Zagros collision: recent advances in its structure and active deformation. *Geol Soc Lond Spec Publ* 330(1):43–64
- Sarwar G, Alizai A (2013) Riding the mobile Karachi arc, Pakistan: understanding tectonic threats. *J Himal Earth Sci* 46(2):9–24
- Scordilis EM (2006) Empirical global relations converting M<sub>S</sub> and m<sub>b</sub> to moment magnitude. *J Seismol* 10(2):225–236
- Shi Y, Bolt BA (1982) The standard error of the magnitude-frequency b-value. *Bull Seismol Soc Am* 72:1677–1687
- Smith WD (1981) The b-value as an earthquake precursor. *Nature* 289(5794):136
- Smith LS (2013) The structure, fluid distribution and earthquake potential of the Makran subduction Zone, Pakistan. Univ. of Southampton, Department of ocean and earth sciences, PhD thesis abstract

- Smith GL, McNeill LC, Wang K, He J, Henstock TJ (2013) Thermal structure and megathrust seismogenic potential of the Makran subduction zone. *Geophys Res Lett* 40:1528–1533
- Szeliga W, Bilham R, Kakar DM, Lodi SH (2012) Interseismic strain accumulation along the western boundary of the Indian subcontinent. *J Geophys Res* 117(1):08404. <https://doi.org/10.1029/2011jb008822>
- Waseem M, Khan MA, Khan S (2019) Seismic sources for southern Pakistan and seismic hazard assessment of Karachi. *Nat Hazards* 99(1):511–536
- Wiemer S (2001) A software package to analyze seismicity: ZMAP. *Seismol Res Lett* 72(3):373–382
- Yeats RS, Lawrence RD, Jamil-ud din S, Khan SH (1979) Surface effects of the 16 March 1978 earthquake, Pakistan–Afghanistan border. In: Farah A, DeJong KA (eds) *Geodynamics of Pakistan*. Geological Survey of Pakistan, Quetta, pp 159–361
- Zinke R, Hollingsworth J, Dolan JF (2014) Surface slip and of fault deformation patterns in the 2013 M 7.7 Balochistan, Pakistan earthquake: implications for controls on the distribution of near-surface coseismic slip. *Geochem Geophys Geosyst* 15:5034–5050



# The use of QLARM to estimate seismic risk in Kirghizstan at the regional and city scales

Philippe Rosset<sup>1</sup> · Stavros Tolis<sup>1</sup> · Max Wyss<sup>1</sup>

Received: 28 November 2019 / Accepted: 23 May 2020 / Published online: 6 June 2020  
© Institute of Geophysics, Polish Academy of Sciences & Polish Academy of Sciences 2020

## Abstract

An analysis of seismic risk using our tool QLARM has been performed for the Batken region including the cities of Aidarken and Kadamjay, 100 km SW of Osh. The damage to residential buildings and induced casualties has been estimated for a set of seismic scenarios of typical and maximum magnitude considering the existing seismicity data. Population and building datasets have been built based on up-to-date information, and for the two cities, satellite photographs and a field survey have been used. A preliminary soil response zonation is proposed using seismic ambient noise analyses. In the investigated region, the probability of damaging earthquakes with  $M > 6$  is judged to be low because the slip accumulation rate along individual faults is only in the range of 0.01–0.3 cm/year. The amplification of seismic waves by soil deposits is estimated to be low; however, the proposed zonation needs to be complemented by additional seismic measurements. The calculations indicate that the combined fatalities of Kadamjay and Aidarken in a hypothetical earthquake of magnitude between 6.0 and 6.6 are fewer than 100.

**Keywords** Earthquake · Risk · Residential buildings · Central Asia · Mitigation

## Introduction

Estimating the damage and human losses due to possible future large earthquakes is of prime importance for mitigation and preparedness. Wyss (2017) pointed out that the number of losses could be largely reduced if loss estimations were used to retrofit critical buildings and promote safety and evacuation measures. The author noted that each US\$1 spent on earthquake mitigation saves \$10 when the disaster strikes.

Kirghizstan is an earthquake-prone country located within the Tien Shan mountain belt. The on-going deformation of the region is associated with the continental collision of the Eurasian and Indian plates which generated several large historic thrust and reverse-faulting earthquakes (Xu et al. 2006). According to Zhang et al. (2004), crustal velocities are around 5 mm/yr across the southern Tien Shan inducing shallow earthquakes. The map of Fig. 1 shows the earthquakes reported and recorded from up to 2014, included

in the dataset of the Central Asia Seismic Risk Initiative, CASRI (Abdrakhmatov 2009).

In the Batken Province, SW of Kirghizstan, where the study was conducted, most of the earthquakes occur in the Pamir region as shown in Fig. 1. Ischuk et al. (2018) report that 75% of the focal mechanisms in this region are thrust faulting with shallow depth ( $> 45$  km).

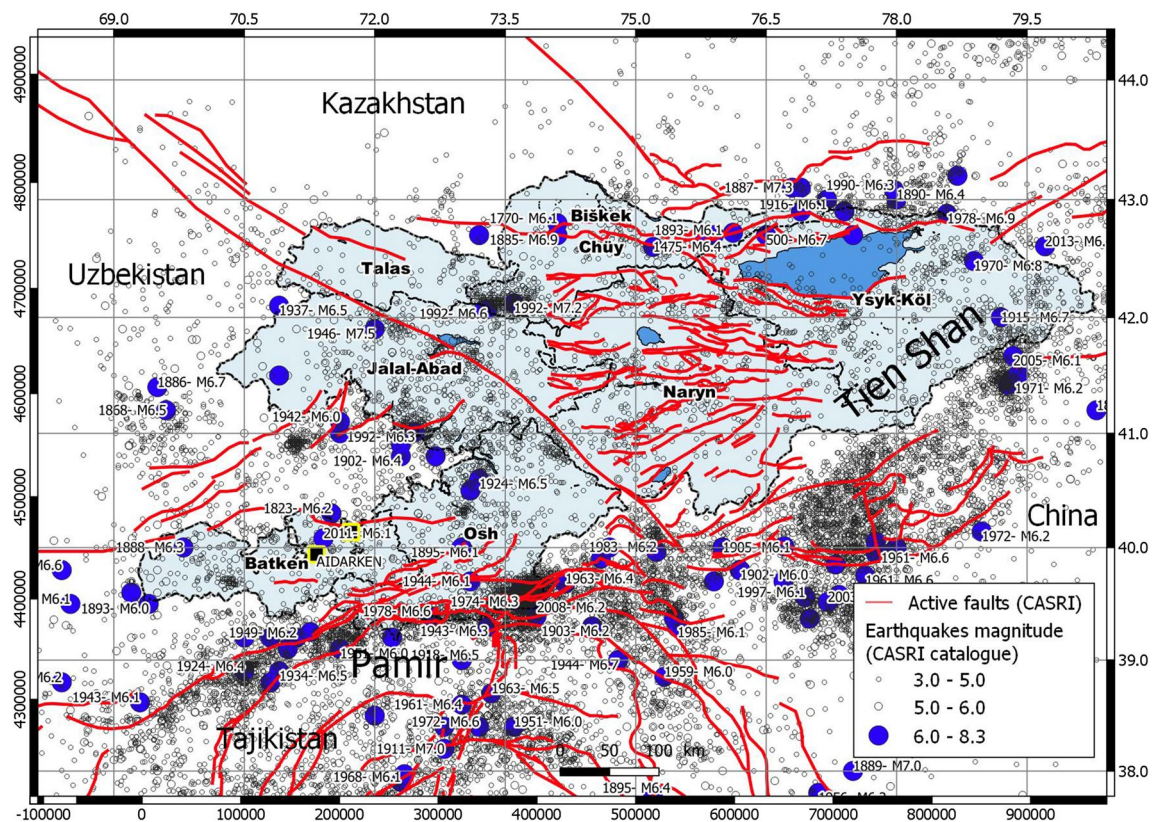
For the last 16 years, we have estimated the extent of earthquake disasters within less than 1 h for about 1200 earthquakes worldwide (Wyss 2014) using our tool QLARM. Loss estimates contain the following information: a map showing the average degree of damage in settlements near the epicenter, a list of the number of people living in areas shaken by each intensity grade of V and larger, the total number of fatalities and the total number of injured.

The database of QLARM includes information for about 1.93 million settlements containing the following parameters. (1) Up-to-date population numbers. (2) Distribution of buildings into EMS-98 classes. (3) Distribution of the population in buildings of these classes. (4) Fragility curves for building resistance to strong ground motions. (5) Occupancy rates for different periods of the day. The program and data sets of QLARM are detailed in several publications (Rosset et al. 2015; Trendafiloski et al. 2009, 2011).

✉ Philippe Rosset  
rossetp@orange.fr

<sup>1</sup> ICES - International Centre for Earth Simulation  
Foundation, Geneva, Switzerland





**Fig. 1** Seismotectonic context of Kirghizstan (shaded in blue) and surrounding countries. Earthquakes listed in the CASRI catalog up to 2014 are located by dots and grouped by magnitude ranges. For earthquakes with magnitude higher than 6 the year and magnitude  $M$  are indicated. Faults are marked with red lines as given by the CASRI

dataset. The Kirghizia border is shown as a black line, and the province borders are named in bold and marked by dashed lines. The two cities, Kadamjay and Aidarken, are indicated with a black square and yellow contour

Several global and regional attenuation relationships of seismic waves are available to calculate peak ground acceleration (PGA) or intensity. For more than 50 important cities, information on soil conditions is provided to consider the effects of site amplification in the ground shaking calculation. In these cases, the settlement is divided into districts for which site amplification parameters are given as well as the population and building parameters described previously (Parvez and Rosset 2014).

QLARM has also been used to estimate losses for likely earthquakes in different regions of the world such as in the Azores (Fontiela et al. 2020), the Himalayas (Wyss and Chamlagain 2019; Wyss et al. 2018b), North India (Wyss et al. 2017), Algeria (Rosset and Wyss 2017), Mexico (Wyss and Zuniga 2016), Southern Sumatra and Central Chile (Wyss 2010) and Central Myanmar (Wyss 2008). The reader will also find in these publications details on QLARM.

The results presented in this paper concern the Batken region (population around 500,000) and most specifically the two cities of Kadamjay and Aidarken with population 13,000 and 11,000, respectively. The project aimed at

compiling data to build a dataset at regional and city scales used in QLARM for assessing the potential human losses and damage to residential building in the cases of various earthquake scenarios. The approach used in this project and the results presented in the paper are detailed in two internal reports (Wyss et al. 2018a, b; Torgoev et al. 2019).

## Construction of the database

An extended compilation of data has been engaged in order to create the hazard and vulnerability models for the Batken oblast and the two investigated cities, Kadamjay and Aidarken.

## Settlement and population dataset

A dataset of 621 settlements in a radius of 200 km around the two cities, including the geographical location, name in Latin and Cyrillic alphabets and population, is the base of our loss models.

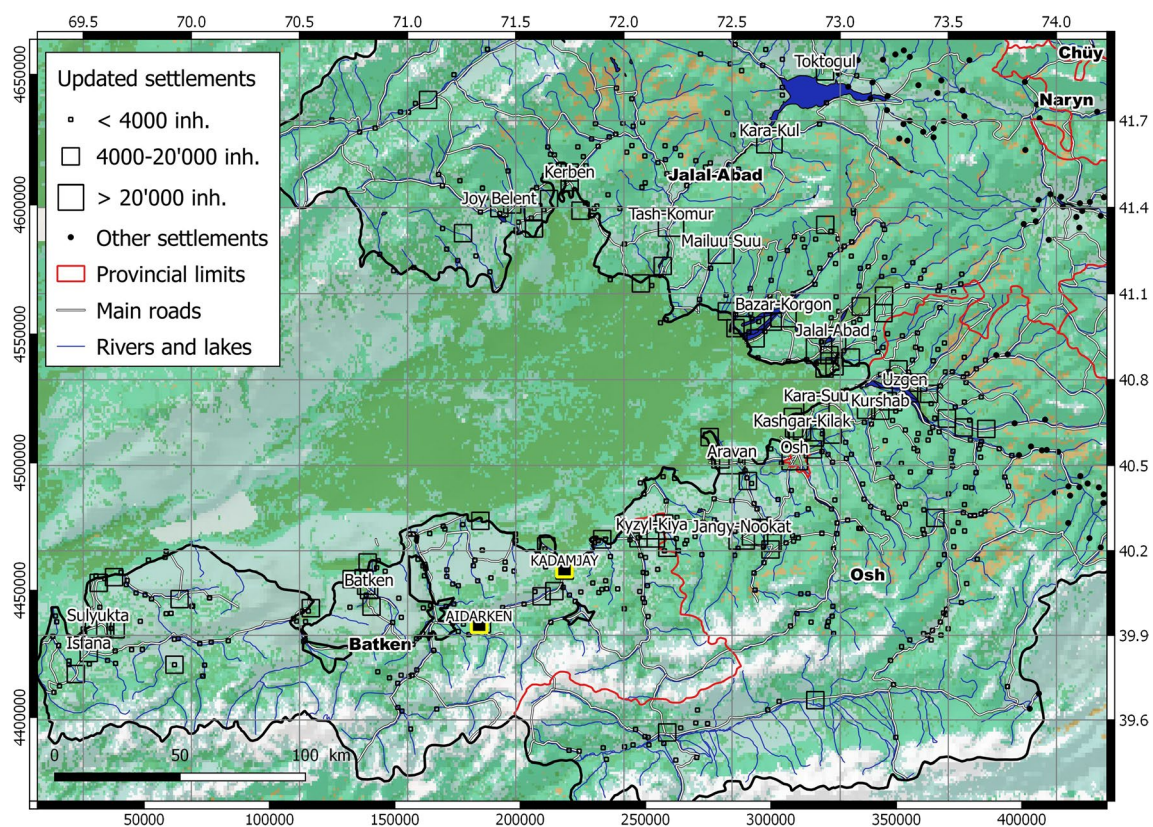
Population data updated for 2017 are derived from (a) official Kyrgyz sources for 200 settlements counting for 67% of the total population, (b) OpenStreetMap (2017) and Global Human Settlement Layer (2016) database for 93 settlements representing 11% of the total population and (c) visual counting of the number of houses in villages from satellite images for the remaining 11% of the population. The total population for the Batken oblast of our dataset is under-estimated by 2% compared to the 2017 official value of 503,500. For the Kadamjay district, our dataset over-estimated the official population count by 5%. These differences are in the range of the expected uncertainty. In total, 85% of the settlements are rural, with fewer than 4000 people per village, 13% middle size and 2% urban settlements with population larger than 20,000 inhabitants. The map of Fig. 2 shows the locations of the updated settlements grouped into these three city sizes for the updated dataset.

### Residential buildings model at regional scale

At the level of Batken oblast, the distribution of buildings by construction types as described in Table 1, is based on

the survey conducted by us on site and the literature (e.g., Lang et al. 2018; Wieland et al. 2015; Tolis et al. 2013; Wyss et al. 2013). Adobe houses (ADO) represent on average two-thirds of residential buildings in the three city sizes, followed by reinforced or confined masonry (RM) and unreinforced fired brick masonry (URM) structures. RM buildings are modeled to contribute 30% of the stock in urban settlements and 17.5% in rural ones. Precast concrete frame structures (RCPC) are mainly found in urban settlements. Photographs of Fig. 3 show the main types of buildings surveyed in the field.

A distribution of residential buildings in terms of the six EMS-98 classes (Gruenthal 1998) of decreasing vulnerability, from A to F, for the three city sizes is derived from the count by building types using the percentages proposed in Table 2. For reinforced concrete moment frame (RC), RM and RCPC buildings, different vulnerability distributions were adopted, depending on the construction date in order to reflect the different construction methods applied during each period considered (1950–1970, 1970–1990 and later). The resulting distributions are given in Fig. 4.



**Fig. 2** Settlements database. Settlements updated with 2017 population are shown with unfilled squares. Black dots are settlements available in the QLARM database outside the 200 km radius zone. Names are indicated for cities with more than 15,000 inhabitants. The Kir-

ghizia border is shown as a black line, and the province borders are named in bold and marked by red lines. The two cities, Kadamjay and Aidarken, are located by a black square with yellow contour

**Table 1** Typical residential building types in Batken oblast

Types	Description
ADO	Adobe block (unbaked dried earth brick), mud mortar, wood roof and floors structures
W	Timber structures
URM1	Unreinforced fired brick masonry, cement mortar, timber flooring structures
URM2	Unreinforced fired brick masonry, cement mortar, (precast) concrete flooring structures
RM	Reinforced or confined masonry structures
RC	Reinforced concrete moment frame with masonry infill wall structures with various level of earthquake-resistant design
RCPC	Precast concrete frame structures

**Fig. 3** Typical buildings in the investigated area. The abbreviation indicated for each building type is explained in the text



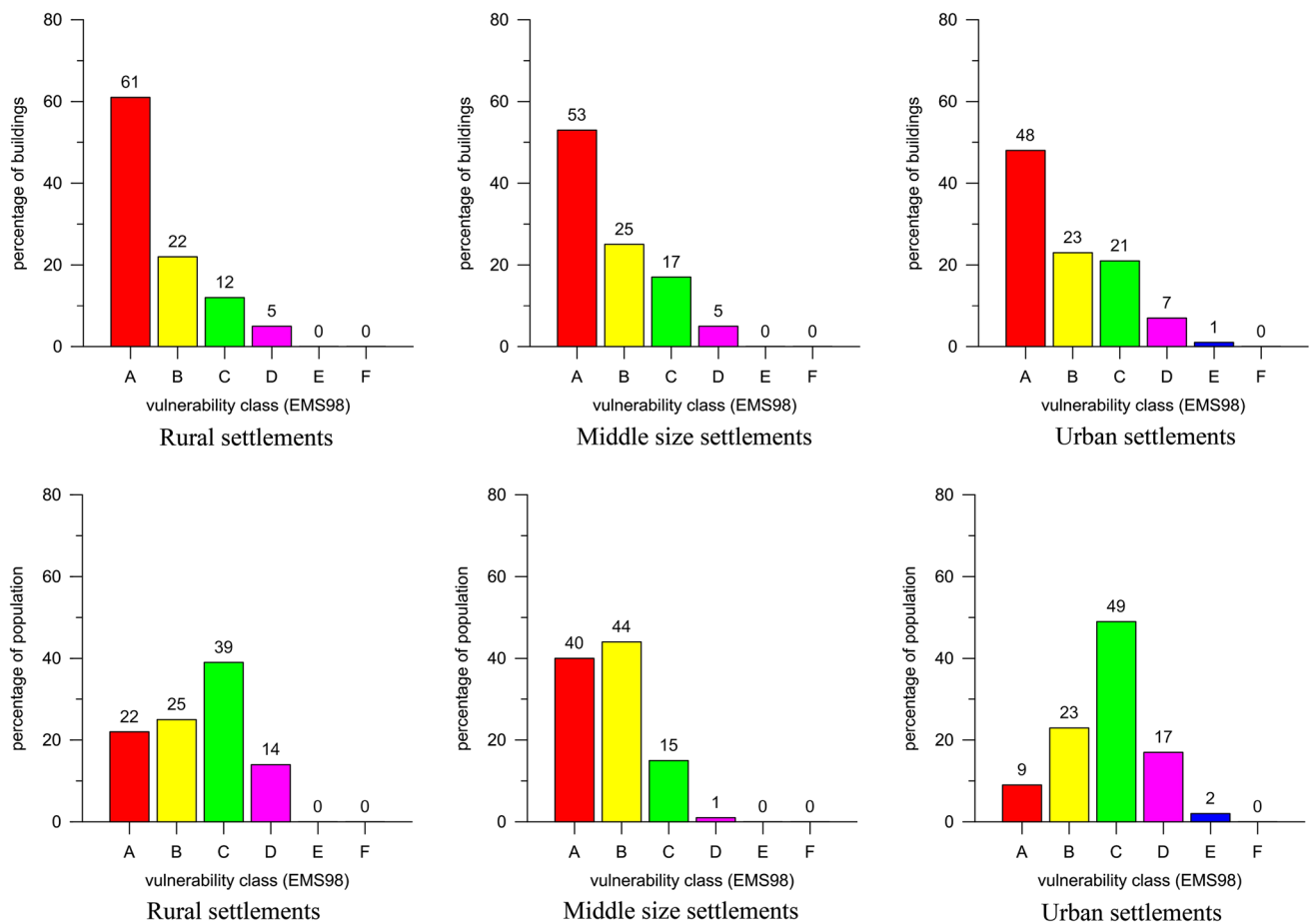
**Table 2** EMS-98 vulnerability distribution for typical construction types for the Batken oblast (in %)

Building types	EMS-98 vulnerability classes					
	A	B	C	D	E	F
ADO	80	20	0	0	0	0
W	0	0	20	60	20	0
URM1	10	80	10	0	0	0
URM2	20	60	20	0	0	0
RM (1950–1990)	0	30	60	10	0	0
RM (1990–)	0	10	50	40	0	0
RC (1950–1970)	0	40	50	10	0	0
RC (1970–1990)	0	20	50	30	0	0
RC (1990–)	0	0	30	50	20	0
RCPC (1950–1990)	20	50	30	0	0	0
RCPC (1990–)	0	0	35	50	15	0

**City models for Aidarken and Kadamjay**

The city models for Aidarken and Kadamjay were developed using the information collected in the field. The

mobile application Kobotoolbox (2019) and a digital model created using the footprints of more than 8100 building as they appear in the online satellite images have been used to document typical residential buildings and



**Fig. 4** Residential vulnerability distributions in terms of buildings (top) and population (bottom) for the three city sizes in Batken oblast

districts in both cities (See the example for Aidarken in Fig. 5).

The response of the soil was investigated using ambient noise records on 82 sites, 39 in Aidarken and 43 in Kadamjay because no other information to estimate the soil conditions locally was available and the number of days in the field was limited. The horizontal to vertical spectral ratio (HVSr) method was used to analyze the records and define zones with similar response frequencies which are often inversely correlated with the thickness of recent deposits (e.g., Lunedei and Malischewsky 2015).

The city of Aidarken is located on a flat EW deposition cone created by alluvial material carried from mountains in the south by the Gauyan river. The alluvial deposits are made of boulder, gravel and sand materials of different levels of compaction. The eastern part of the town is built on loess deposits (Neogene period) which are progressively eroded by seasonal streams.

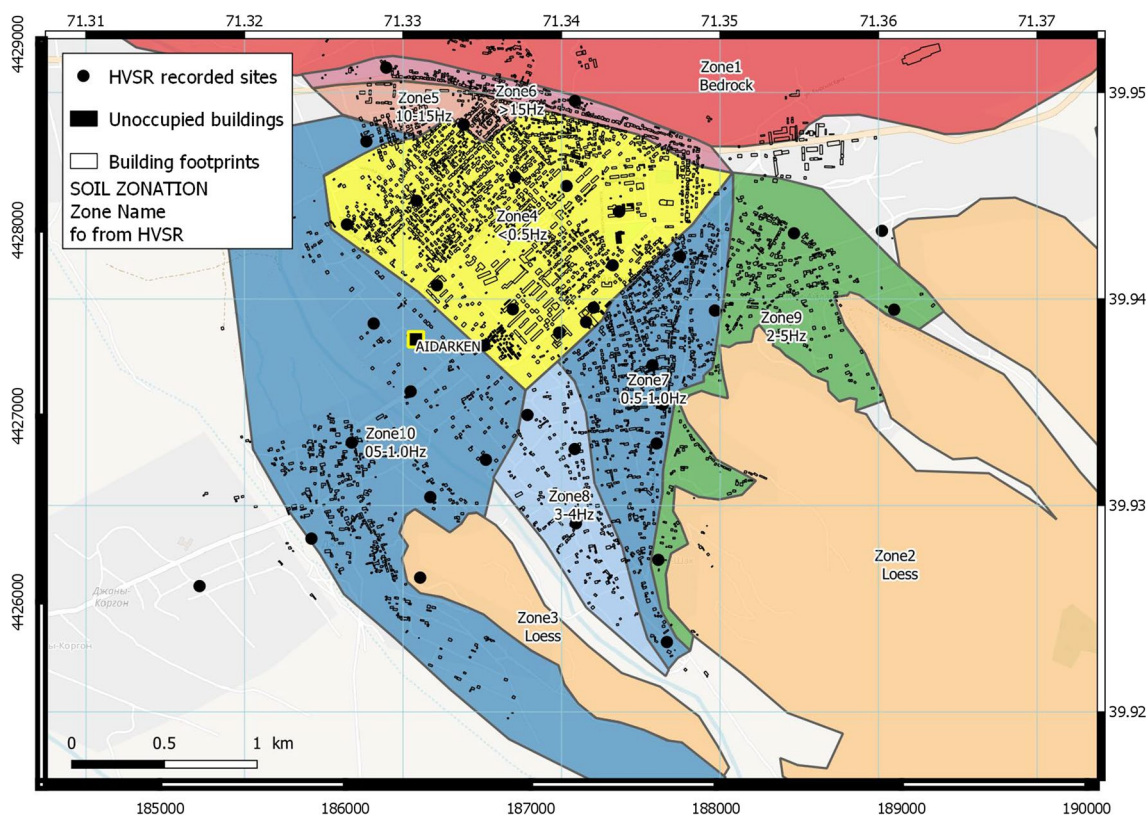
High mountains surround the valley to the south and north. The distribution of predominant frequency  $f_0$  from the HVSr analysis is well correlated with the history of alluvial

depositions coming from the southern mountain streams via the SE part of the city (Fig. 5).

The northern part of the city exhibits a high-frequency peak in the ambient noise records indicating a thin layer of soft sediment overlying hard rock constituting the mining site. Sites in the center of the city show low-frequency peaks (0.3–1.5 Hz) indicating a thick layer of deposits. In the eastern part of the city, higher-frequency peaks reveal more recent thin deposits from the river or from loess erosion.

For Aidarken, the model is a normal settlement because the buildings are mostly single-family houses except in a quarter with a tens of RM five-floors buildings (Table 3) and the site amplification of the coarse alluvial deposits in the built areas is suspected to be low.

In Kadamjay, the soil response is influenced by the thickness of the terraces in the eastern and western part of the city. The map of Fig. 6 shows the proposed zonation in terms of predominant response frequency  $f_0$  in the HVSr spectra. The built environment in conjunction with site response characteristics allowed the division of Kadamjay in four zones delimited in bold lines in Fig. 6:



**Fig. 5** City model for Aidarken. The map shows the footprints of buildings with unoccupied ones filled in black. Black dots locate the sites where ambient noise has been recorded. The zonation in terms of predominant frequency  $f_0$  (in Hz) is visible as background colors

**Table 3** Building distribution by construction types in Kadamjay and Aidarken

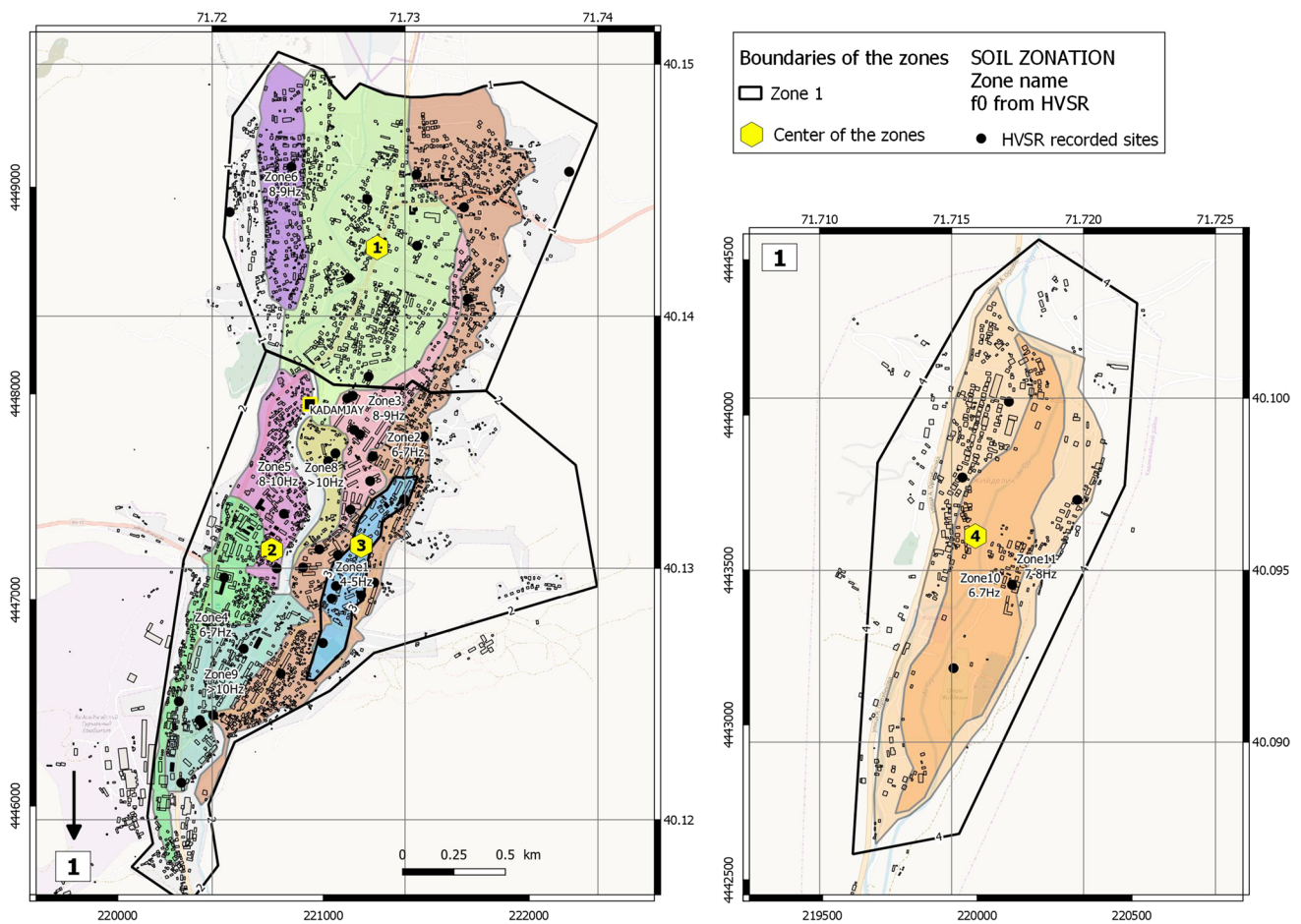
Kadamjay	Building distribution by construction type (%)						
	ADO	W	URM1	URM2	RM	RC	RCPC
Zone 1	68.0	1.0	6.0	0.0	23.0	2.0	0.0
Zone 2	64.5	0.0	5.5	5.0	17.0	7.5	0.5
Zone 3	38.5	0.0	13.0	17.0	19.0	8.5	4.0
Zone 4	85.0	1.0	5.0	4.0	5.0	0.0	0.0
Aidarken	65.0	1.5	5.0	2.0	22.0	4.0	0.5

- *Zone 1* in the northern part where  $f_0$  is higher than 6 Hz; the built environment consists mainly of single-family houses with 1 or 2 floors.
- *Zone 2* in the southern part where  $f_0$  varies between 6 and 10 Hz; is covered with a mixture of commercial and residential buildings of fewer than 4 floors.
- *Zone 3* is located inside zone 2 in the eastern section of Kadamjay, where most of the buildings are classified as URM, RM and RCPC with 4–5 floors. The estimated response frequency  $f_0$  varies between 4 and 5 Hz, values which are close to the dominant frequency of the buildings in this zone. In order to consider possible damaging resonance effects between soil and structure responses, the calculated intensities were increased adding 0.2 units;

this value is empirical, arrived by expert judgment, and independent of the level of shaking; dynamic structural analyses should verify in a future stage the empirical value used.

- *Zone 4* comprises Jydelik, a district south of Kadamjay, with  $f_0$  between 6 and 8 Hz; most of the buildings are single-family houses with 1 or 2 floors.

The city models of Kadamjay and Aidarken are based on the scale of individual blocks thanks to the available footprint data used during the field survey. The most typical buildings in the different zones of the cities were visited in order to count and distribute the footprints dataset by construction types. Buildings with a footprint surface lower



**Fig. 6** Division of the city of Kadamjay into zones. Four zones are delimited by bold lines in Kadamjay, three in the town itself (left map) and one in the south (right map marked 1). They are represented in the QLARM dataset by a X, Y coordinates (yellow hexagons). The

maps show the footprint of buildings with unoccupied buildings filled in black. Black dots locate the sites where ambient noise has been recorded. The zonation in terms of predominant frequency  $f_0$  is visible as background colors

than 45 m<sup>2</sup> were counted as unoccupied. Table 3 lists the distribution by construction types for the zones identified in both cities.

The building distribution in terms of EMS-98 presented in Table 4 for Kadamjay and Aidarken is based on the detailed city models described in Tables 2 and 3. In addition, the number of apartments in each building and the average occupancy rates were counted in order to define a distribution of the population living in buildings of the EMS-98 vulnerability classes (Table 4).

### Risk estimates at the city scale for Kadamjay and Aidarken

#### Calculation of the ground shaking in terms of intensity

The amount of ground shaking due to an earthquake can be estimated as intensity (I), peak ground accelerations (PGA) or other parameters. For scenario calculations, QLARM uses most often intensities, because these values can be compared to observed intensities. This serves as verification of the algorithms and data in QLARM. Among the numerous equations to calculate intensities, the average equation used here has been proposed by Shebalin (1968, 1985) and is described by:

$$I = C1 * M - C2 * \log(\text{sqrt}(R^2 + h^2)) + C3 \quad (1)$$

**Table 4** Distribution of buildings and population in terms of EMS-98 classes for Kadamjay and Aidarken

Kadamjay	Building distribution in EMS-98 classes (%)					
	A	B	C	D	E	F
Zone 1	59.9	21.9	15.1	3.0	0.0	0.0
Zone 2	57.7	27.2	12.5	2.7	0.0	0.0
Zone 3	38.3	44.7	17.0	0.0	0.0	0.0
Zone 4	70.0	25.0	5.0	0.0	0.0	0.0
Aidarken	59.8	26.7	12.5	1.0	0.0	0.0
Kadamjay	Population distribution in EMS-98 classes (%)					
	A	B	C	D	E	F
Zone 1	48.2	27.3	20.2	4.3	0.0	0.0
Zone 2	29.2	52.4	15.4	3.0	0.0	0.0
Zone 3	13.9	80.4	5.7	0.0	0.0	0.0
Zone 4	70.0	25.0	5.0	0.0	0.0	0.0
Aidarken	34.4	53.0	11.9	0.8	0.0	0.0

where  $M$  is the magnitude,  $R$  is the closest distance to the rupture plane in km,  $h$  is the depth of the hypocenter in km, and  $C_1$ ,  $C_2$  and  $C_3$  are constants. This equation developed specifically for former USSR countries is used because it can easily be adjusted to local conditions by changing the constants. For this region, the constants were not adapted because there was no need as explained in Sect. “Validation of the QLARM model for Central Asia”; the standard average values 1.5, 4.5 and 3.5 were used for  $C_1$ ,  $C_2$  and  $C_3$ , respectively.

### Validation of the QLARM model for Central Asia

Calculated intensities were compared to those shown on maps for three earthquakes reported by Kalmetieva et al. (2009). These earthquakes were the 1911 M8.2 Kebin, the

1946 M7.4 Chatkal and the 1992 M7.2 Suusamyр earthquakes. These events were located in Central Asia and provided an approximate validation of QLARM estimates of intensities when published magnitudes and locations were used, along with reasonable assumptions of shallow depths and average attenuation values. Although these comparisons cannot be considered a rigorous validation for QLARM estimates of intensities in the study area, they provided an approximate verification that QLARM calculates realistic values for Intensity in Central Asia.

The comparison of the estimates of casualties calculated by us with the reported ones allows to gauge to what extent QLARM furnishes reliable result, given only hypocenter and  $M$ . Table 5 lists all earthquakes large enough for a near-real-time alert issued by the QLARM team since the fall of 2002 and which occurred in Kyrgyzstan, Tajikistan and

**Table 5** Comparison between observed (obs.) and our estimated (estim.) casualties for important earthquakes in Central Asia

Year	Month	Day	Lon. (deg)	Lat. (deg)	Depth (km)	$M$	Fatalities (estim.)		Fat (obs.)	Injured (estim.)		Inj (obs.)	Country
							Min	Max		Min	Max		
2004	11	17	71.82	39.19	43	5.7	0	0	0	1	7	–	Tajikistan
2005	2	25	72.70	38.1	107	5.9	0	0	0	0	0	–	Tajikistan
2007	1	8	70.31	39.82	18	5.8	0	3	0	16	146	–	Kyrgyzstan
2008	10	5	73.82	39.53	27	6.7	3	85	74	42	439	140	Tajikistan
2011	1	24	72.85	38.40	102	6.0	0	0	0	0	0	–	Tajikistan
2011	7	19	71.41	40.08	20	6.1	2	50	14	99	935	86	Tajikistan
2013	5	26	67.31	39.96	18	5.7	1	22	0	68	598	–	Uzbekistan
2015	11	17	73.26	40.35	15	5.5	0	1	0	2	82	–	Kyrgyzstan
2015	12	7	72.78	38.21	22	7.2	38	411	2	174	1411	>100	Tajikistan
2016	6	26	73.34	39.48	13	6.4	1	17	2	12	116	14	Tajikistan
2017	5	3	71.44	39.49	11	6.0	0	6	0	5	47	0	Tajikistan

Uzbekistan. The list contains events from these three countries because the building types and the tectonic setting are similar and one needs more cases of comparison than just the two that occurred in Kyrgyzstan.

Table 5 lists the updated latitudes, longitudes, depths and  $M$  after each earthquake and the reported fatalities and injuries as given by the significant earthquakes database of NOAA (2019).

The range of the casualty estimates is large because numerous uncertainties exist (Wyss and Rosset 2013). In cases of small numbers of fatalities, small differences are judged as an adequate result because a single building collapsing or not can change the total numbers of fatalities on the order of the total estimate. Based on the comparison in Table 5, the intensity of shaking and the numbers of casualties are calculated approximately correctly for significant earthquakes in Central Asia, especially in the eastern part of Kyrgyzstan where the database was updated.

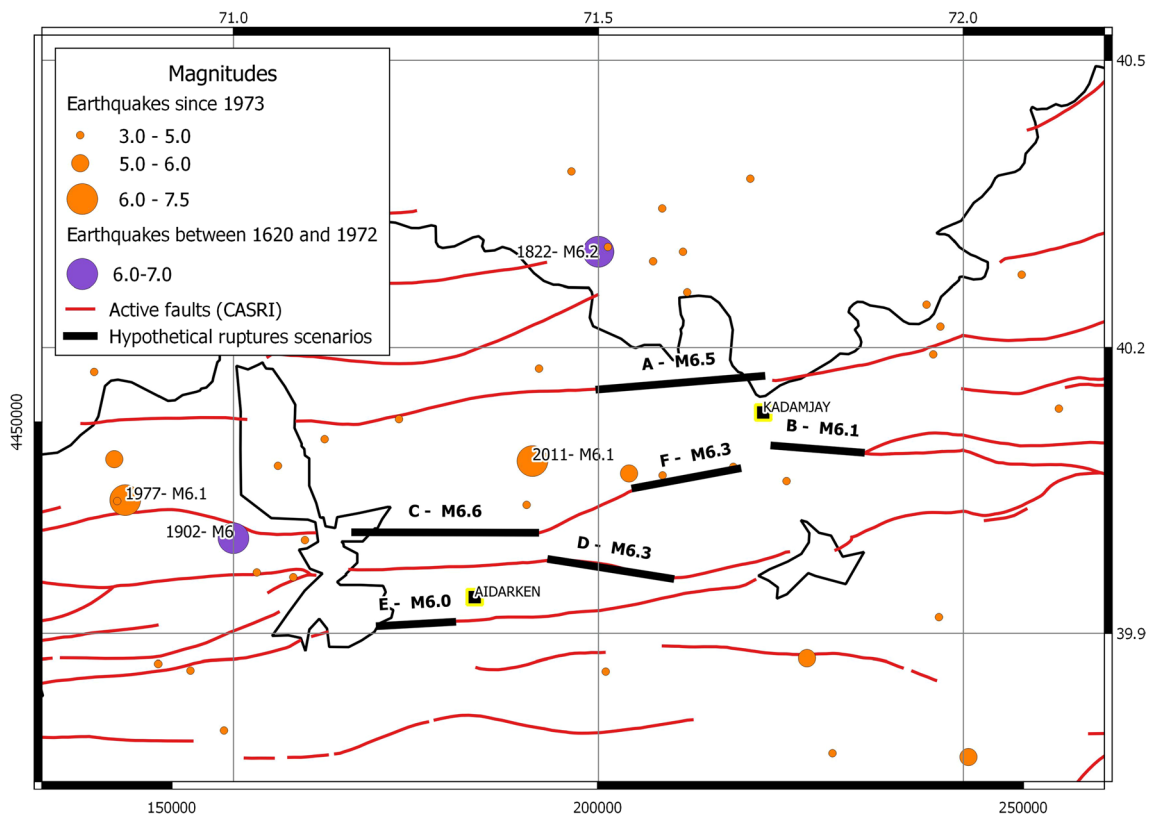
### Definition of the scenarios

There exist numerous maps with large scales showing myriads of active faults in Central Asia (e.g., Abdrakhmatov et al.

2003; Bindi et al. 2012; Ischuk et al. 2018). Most of the faults on these maps are irrelevant to damage in the study area because earthquakes generated along them are too far away. Only those faults matter, which are nearer to our locations of concern than about 40 km.

There are two sources for faults available in the study area. The collection of faults contained in the dataset of CASRI (Abdrakhmatov 2009) has been mapped by various geologists, but no references to these authors are given. The other source is a collection of maps placed on the web by Mohadjer et al. (2016) where the authors are cited. When comparing these maps, one sees that many faults are the same in both datasets, but some are different. The number of faults in the CASRI dataset near the two towns of interest is far greater than that by Mohadjer et al. (2016). Because mapped faults cannot be ignored, even if they should be inaccurate, the CASRI faults are considered as possible sources for earthquakes strongly affecting the study region (Fig. 7).

Two types of scenarios,  $M_{max}$  and  $M_{typical}$ , are proposed based on their magnitude  $M$ ;  $M_{max}$  defines the largest magnitude earthquake that can be expected in a given region.  $M_{typical}$  is defined as the  $M$  assigned to several earthquakes



**Fig. 7** Map of the quaternary faults (red lines) near Kadamjay and Aidarken from the CASRI dataset. Recorded and historic earthquakes are located by red and blue dots, respectively. Selected fault lines for

likely scenarios are marked in black with the name of the scenario and its magnitude  $M$  as listed in Table 5. Aidarken and Kadamjay are located by red and blue dots, respectively. Selected fault lines for



that happened during the period of the seismicity catalog within 100 km of the investigated point and which have caused fatalities. Their magnitudes ranged from 6.2 to 6.3. The hypothetical earthquakes are proposed by selecting fault segments for which end points of a possible rupture can be identified. For defining scenario ruptures, the information available is locations, lengths, azimuth of the strike of the fault and changes of strike. We assume that the faults shown in the map from CASRI are all active. We also assume that changes of strike in faults (kinks) may be the end of a rupture and that the relationships between  $M(\text{length})$  and  $M(\text{area})$  for thrusts given by Wells and Coppersmith (1994) are applicable. The depth is fixed to 16 km which corresponds to the average value of calculated depths in the CASRI catalog for this region. The six scenarios selected, in the range of M6.2–M6.6 (as shown in Fig. 6 and listed in Table 6), are the following:

**Scenario A** The earthquake of M6.1 in 2011 was probably located on the fault north of the estimated epicenter. Assuming the fault had the standard length of an M6.1 earthquake and the epicenter was located at the center of the break, we assign the probable eastern end of this rupture as 71.5E/40.15N. As in every earthquake, stress is transferred to the neighboring sections of the fault in question. Thus, the probability of rupture for the fault segment adjacent and east of this M6.1 earthquake in 2011 are increased. Therefore, we assign as scenario A the fault section from this point to the interruption of this fault on the geologic map (Fig. 7). The rupture ends, the epicenter and the resulting length of 25 km are defined by these assumptions. Further assuming a width of rupture of 10 km the magnitude is estimated as M6.5.

**Scenario B** is an example of a Mtypical earthquake to occur near Kadamjay. The fault segment ends in the west, and in the east bifurcates, limiting its length to 10 km. This results in an M6.2 earthquake for which a width  $W$  of 10 km is assumed.

**Scenario C** Assuming that the 1902 M6 earthquake occurred on the fault mapped just north of its epicenter,

a second fault segment is identified in which stress was recently increased and that ends close to Aidarken. With a length of 20 km and a width of 12 km, its magnitude is estimated as M6.6.

**Scenario D** A fault segment east of Aidarken offers itself as having fairly clearly defined ends that lead to  $L = 14$  km. With an assumed  $W = 10$  km, the magnitude is estimated to 6.3.

**Scenario E** The western end of the fault mapped just south of Aidarken may be capable of rupture in an M6.2 earthquake.

**Scenario F** In the central part of the map of Fig. 6, a fault is located that had may have generated minor seismic activity recently: A possible rupture is defined in its eastern half. With a rupture length of 12 km, one expect that this segment is capable of an M6.3 earthquake.

## Discussion

In QLARM, calculated building damage for each individual settlement is divided into six levels. The maps of Fig. 8 show the calculated damage by degrees (left map) and mean damage (right map) in each settlement around the fault (black dotted line) for the M6.6 scenario C. The damage, divided in six levels, is represented on a pie chart, each color representing the percentage of a certain degree of damage from no damage to collapse. The mean damage grade  $M_d$ , which is the combination of the damage by degree and the distribution of population by vulnerability classes, is given by a colored dot and again divided in six levels, from no damage (0) to complete (5).

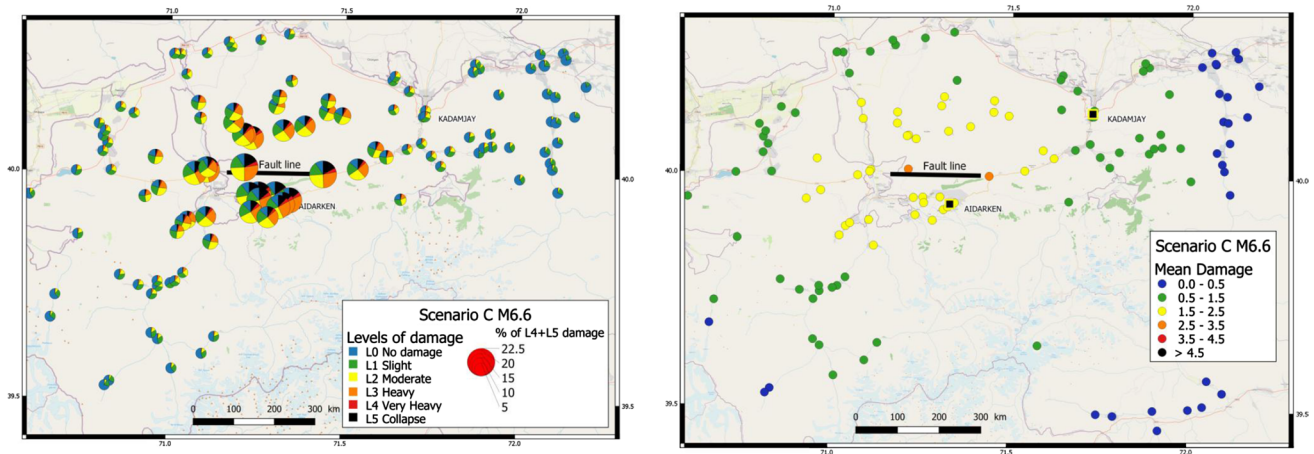
At the scale of Kadamjay and Aidarken, similar results are proposed at a better resolution when the city is divided in zones like in Kadamjay. The map of Fig. 9 shows the calculated damage for the same M6.6 scenario C.

Fatalities that may result in Kyrgyzstan and for the zones of Kadamjay and Aidarken separately are based on the proposed six reasonable, but hypothetical scenarios (Table 6). The estimated numbers of fatalities carry

**Table 6** Proposed scenario parameters for six hypothetical earthquakes in the study area

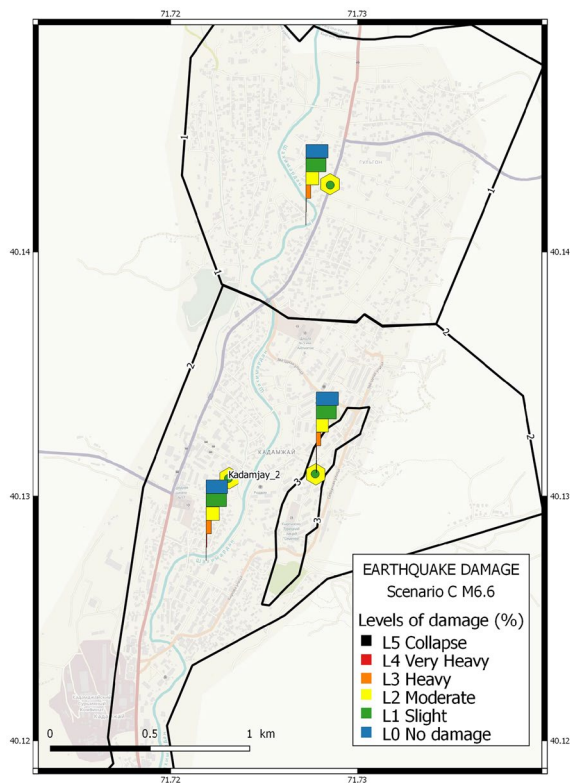
ID	Lon1 (deg)	Lat1 (deg)	Lon2 (deg)	Lat2 (deg)	Lon_epi (deg)	Lat_epi (deg)	Lenght (km)	Width (km)	Depth (km)	Mw
A	71.5	40.15	71.72	40.17	71.61	40.16	20	10	16	6.5
B	71.74	40.10	71.86	40.09	71.80	40.10	10	10	16	6.1
C	71.17	40.00	71.41	40.00	71.29	40.00	20	12	16	6.6
D	71.43	39.98	71.6	39.96	71.52	39.97	14	10	16	6.3
E	71.2	39.91	71.3	39.91	71.25	39.91	8.5	10	16	6.0
F	71.55	40.05	71.69	40.07	71.62	40.06	12	10	16	6.3

The two ends of the assumed ruptures are given as Lon1/Lat1 and Lon2/Lat2, respectively, and the location of the epicenter Lon\_epi/Lat\_epi is also listed. The assumed length ( $L$ ) and width ( $W$ ) result on average in the listed  $M$ , using the relationship by Wells and Coppersmith (1994)



**Fig. 8** Calculated building damage at regional scale in the case of scenario C. Each settlement around the fault line (black thick line) is located by a pie chart or a dot. (Left map) The surface of the pie chart

is proportional to the percentage of heavy damage and collapse. The color corresponds to the level of damage. (Right map) The color of each dot is the mean damage grade as defined in the text



**Fig. 9** Calculated building damage for Kadamjay in the case of scenario C. Each zone with different building distribution and soil conditions of the city model is delimited by bold lines and located by a dot with the color corresponding to the mean damage grade (here green as slight). Vertical histograms show the percentage of damage by degree from L0 to L5. The sum of damage is 100%

uncertainties larger than the averages listed in Table 7 for each scenario. These fatality estimates are calculated for nighttimes, which means maximum occupancy of dwellings, and therefore worst case scenarios result.

Zero fatalities are possible as a minimum in all earthquake scenarios for both cities. On the other hand, the numbers of fatalities could also be factors of 2–3 larger than the values given in Table 7. The sources of uncertainties are multiple as discussed in Wyss (2014). A significant difference in casualties can result from nighttime (high occupancy rate) to daytime earthquakes. Also a collapse of a single apartment building can more than double the numbers of casualties calculated. In addition, there are several earthquake source parameters that can vary and the soil conditions can modify the wave amplitudes. Thus the casualties are calculated with large uncertainties.

The fatalities due to the hypothetical scenarios with identification A–F of earthquakes in Table 5 are listed in Table 7, and the numbers of injured are given in Table 8 for nighttime occurrences. All numbers of estimated casualties are approximate. The numbers of fatalities have at least been verified, but the numbers of injured are unsure because the degree of “injuries” is not defined in reports of casualties in earthquakes. Here, “injured” means that a person needs to be admitted to a health facility.

The numbers of casualties in Tables 7 and 8 are rounded to the nearest integer. Two exceptions are made for very small numbers where value 1 is given instead of 0 when the average was between 1 and 2, and a value of 5 is given when the average was between 2 and 5.

In the present cases, uncertainties, such as population numbers, building type variations, varying soil conditions and especially unknown earthquake source parameters, introduce many unknowns. Therefore, the casualty

**Table 7** Fatalities due to assumed scenario earthquakes (Sc.) during nighttime

Sc.	Kadamjay					Aidarken		Kyrgyzstan		All countries		
	Zone 1	Zone 2	Zone 3	Zone 4	Total		Min	Max	Min	Max	Min	Max
	Average	Average	Average	Average	Min	Max						
A	40	60	20	10	90	180	0	5	150	340	150	2190
B	10	20	5	5	20	60	0	0	30	90	10	300
C	1	1	1	1	0	10	80	150	200	430	180	1980
D	1	5	1	1	0	10	20	60	40	110	10	280
E	0	0	0	0	0	0	10	30	20	50	10	250
F	10	20	5	5	20	60	0	5	40	120	10	310

**Table 8** Injured due to assumed scenario earthquakes (sc.) during nighttime

Sc.	Kadamjay					Aidarken		Kyrgyzstan		All countries		
	Zone 1	Zone 2	Zone 3	Zone 4	Total		Min	Max	Min	Max	Min	Max
	Average	Average	Average	Average	Min	Max						
A	140	270	60	30	360	650	20	40	940	1900	1050	9850
B	40	110	20	20	120	270	1	5	280	620	210	2140
C	10	20	5	5	20	50	310	600	1120	2250	1050	8800
D	10	20	5	5	20	60	110	240	340	750	220	2160
E	1	1	0	0	1	5	60	140	160	370	120	1320
F	40	100	20	20	110	250	10	30	360	780	240	2370

estimates in Tables 7 and 8 are only order of magnitude estimates.

The estimated number of injured is typically about three times that of the fatalities. This is true for cases with large numbers of casualties. However, when the fatality estimate is near zero, the number of injured may still be substantial, that is 100 or more injured may be reported when there are no fatalities. The numbers of injured are given in Table 8. An “injured” may be counted if the person is admitted as a patient in a hospital, but also if the person is treated and released as an outpatient.

The main results are that (1) in the towns of Kadamjay and Aidarken in all of the proposed scenarios the expected fatalities are below 100, and (2) that only the two largest magnitude disasters generate moderate to serious numbers of fatalities (Table 7).

The numbers of injured people, however, are larger than 100 in all cases, and they could number several thousand for the larger events (Table 8). If the medium size and large earthquake scenarios should happen, the country of Kyrgyzstan would face a serious problem of taking care of so many patients.

We have made a great deal of progress toward understanding seismic risk in Kadamjay, Aidarken and the surrounding area. However, the resources were limited and hence the results are not as complete as we would wish. To support the loss estimates for these hypothetical

earthquakes, a detailed and careful review of local building types and an upgrade of the regional population distribution have been performed. As a result, the deterministic loss estimates due to earthquakes have become more reliable than most estimates of casualties and numbers of injured. Nevertheless, the accuracy of these case estimates should not be overestimated. A seismic measurements campaign should be carried out to assess the shear-wave velocity  $V_s$  and thickness of the different soil layers in the eastern terrace of Kadamjay and in selected zones of Aidarken. The building vulnerability database for residential and commercial buildings should be completed by launching local campaigns of crowd sourcing.

**Acknowledgements** The project was funded by Médecins Sans Frontières (MSF) office in Geneva, Switzerland. It benefits from the support of Philippe Calain and the logistic of the MSF team in Kirgizstan and in Geneva. The project was under the auspice of the Kyrgyz Minister of Emergency. We thank the reviewers for helpful comments.

## References

- Abdrakhmatov K (2009) ISTC Project No. KR 1176, Establishment of the Central Asia Seismic Risk Initiative (CASRI). Technical report on the work performed from: 02.01.2006 to 04.30.2009, Institute of Seismology, National Academia of Sciences, Kyrgyz Republic

- Abdrakhmatov K, Havenith H-B, Delvaux D, Jongmans D, Trefois P (2003) Probabilistic PGA and Arias intensity maps of Kyrgyzstan (Central Asia). *J. Seismolog.* 7:203–220
- Bindi D, Abdrakhmatov K, Parolai S, Mucciarelli M, Gruenthal G, Ischuk A, Mikhailova N, Zschau J (2012) Seismic hazard assessment in Central Asia: outcomes from a site approach. *Soil Dyn Earthq Eng* 37:84–91
- Fontiela J, Rosset P, Wyss M, Bezzeghoud M, Rodrigues F (2020) Human losses and damage expected in future earthquakes in Faial Island—Azores. *Pure Appl Geophys* 177:1831–1844
- Global Human Settlement Layer (2016) Corbane C, Florczyk A, Pesaresi M, Politis P, Syrris V (2018) GHS built-up grid, derived from Landsat, multitemporal (1975–1990–2000–2014), R2018A. European Commission, Joint Research Centre (JRC) <https://doi.org/10.2905/jrc-ghsl-10007>
- Gruenthal G (1998) European macroseismic scale. Conseil de l'Europe, Luxembourg
- Ischuk A, Bjerrum LW, Kamchybekov M, Abdrakhmatov K, Lindholm C (2018) Probabilistic seismic hazard assessment for the area of Kyrgyzstan, Tajikistan, and Eastern Uzbekistan, Central Asia. *Bull Seis Soc Am* 108(1):130–144
- Kalmetieva ZA, Mikolaichuk AV, Moldobekov BD, Meleshko AV, Jantaev MM, Zubovich AV (2009) Atlas of earthquakes in Kyrgyzstan. Technical report UNISDR, ISBN 978-9967-25-829-7
- Kobotoolbox (2019) <https://www.kobotoolbox.org>
- Lang DH, Kumar A, Sulaymanov S, Meslem A (2018) Building typology classification and earthquake vulnerability scale of Central and South Asian building stock. *J Build Eng* 15:261–277
- Lunedei E, Malischewsky P (2015) A review and some new issues on the theory of the H/V technique for ambient vibrations. In: Ansal A (ed) Perspectives on European earthquake engineering and seismology. Geotechnical, geological and earthquake engineering, vol 39. Springer, Cham, pp 371–394
- Mohadjer S, Ehlers TA, Bendick R, Stübner K, Strube T (2016) A Quaternary fault database for central Asia. *Nat Hazards Earth Syst Sci* 16:529–542
- NOAA (2019). National Geophysical Data Center/World Data Service (NGDC/WDS): Significant Earthquake Database. National Geophysical Data Center, NOAA. <https://doi.org/10.7289/V5TD9V7K>
- OpenStreetMap (2017) <https://www.openstreetmap.org>
- Parvez I, Rosset P (2014) The role of microzonation in estimating earthquake risk. In: Earthquake, hazard, risk and disaster, Elsevier's hazards and disaster series, pp 273–308
- Rosset P, Wyss M (2017) Seismic loss assessment in Algeria using the tool QLARM. *Civil Eng Res J.* <https://doi.org/10.19080/CERJ.2017.02.555583>
- Rosset P, Bonjour C, Wyss M (2015) QLARM, un outil d'aide à la gestion du risque sismique à échelle variable. In: Leone F, Vinet F (eds) Plan de sauvegarde et outils de gestion de crise. Presses Universitaires de la Méditerranée, Collection Géorisques, Montpellier, pp 91–98
- Shebalin NV (1968) Methods of engineering seismic data application for seismic zoning. In: Medvedev SV (ed) Seismic zoning of the USSR. Science, Moscow, pp 95–111
- Shebalin NV (1985) Regularities of the natural disasters (in Russian). *Nauki o zemle, Znanie* 11:48
- Tolis S, Rosset P, Wyss M (2013) Detailed building stock at regional scale in three size categories of settlements for 18 countries worldwide, Geneva, Switzerland. UNISDR report, 87 pages and appendix, <https://www.unisdr.org/we/inform/publications/49798>
- Torgoev I, Havenith HB, Wyss M., Rosset P, Tolis S (2019) Оценки сейсмической опасности в Баткенской области и сопутствующих рисков в Кадамжае и Айдаркене// Мониторинг, прогнозирование опасных процессов и явлений на территории Кыргызской Республики (Изд. 16-е с изм. и доп.), Бишкек: МЧС КР, pp 688–709 (in Russian)
- Trendafiloski G, Wyss M, Rosset P, Marmureanu G (2009) Constructing city models to estimate losses due to earthquakes worldwide: application to Bucharest Romania. *Earthq Spectra* 25(3):665–685
- Trendafiloski G, Wyss M, Rosset P (2011) Loss estimation module in the second generation software QLARM. In: Spence R, So E, Scawthorn C (eds) Human casualties in earthquakes: progress in modeling and mitigation. Springer, Berlin, pp 381–391
- Wells DL, Coppersmith KJ (1994) New empirical relationships among magnitude, rupture length, rupture width, rupture area and surface displacement. *Bull Seismol Soc Am* 84(4):974–1002
- Wieland M, Pittore M, Parolai S, Begaliev U, Yasunov P, Tyagunov S, Moldobekov B, Saidiy S, Ilyasov I, Abakanov T (2015) A multi-scale exposure model for seismic risk assessment in central Asia. *Seismol Res Lett* 86(1):210–222
- Wyss M (2008) Estimated human losses in future earthquakes in central Myanmar. *Seismol Res Lett* 79(4):504–509
- Wyss M (2010) Predicting the human losses implied by predictions of earthquakes: Southern Sumatra and Central Chile. In: Savage MK, Rhoades DA, Smith EGC, Gerstenberger MC, Vere-Jones D (eds) Seismogenesis and earthquake forecasting: The Frank Evison Volume II. Pageoph Topical Volumes. Springer, Basel. <https://doi.org/10.1007/s00024-010-0090-4>
- Wyss M (2014) Ten years of real-time earthquake loss alerts. In: Wyss M (ed) Earthquake Hazard, Risk, and Disasters. Elsevier, Waltham, pp 143–165
- Wyss M (2017) Reported estimated quake death tolls to save lives. *Nature* 545(7653):151–153
- Wyss M, Chamlagain D (2019) Estimated casualties in possible future earthquakes south and west of the M7.8 Gorkha earthquake of 2015. *Acta Geophys* 67:423–429
- Wyss M, Rosset Ph (2013) Mapping seismic risk: the current crisis. *Nat Hazard* 68(1):49–52
- Wyss M, Zuniga R (2016) Estimated casualties in a possible great earthquake along the Pacific coast of Mexico. *Bull Seismol Soc Am* 106(4):1867–1874
- Wyss M, Tolis S, Rosset P, Pacchiani F (2013) Approximate Model for Worldwide Building Stock in Three Size Categories of Settlements, Geneva, Switzerland. UNISDR report, 34 pages and appendix, <https://www.preventionweb.net/english/hyogo/gar/2013/en/bgdocs/WAPMERR,%202012.pdf>
- Wyss M, Gupta S, Rosset P (2017) Casualty estimate in two up-dip complementary Himalayan earthquakes. *Seismol Res Lett* 86(6):1508–1515
- Wyss M, Rosset P, Tolis S, Havenith HB, Torgoev I, Speiser M (2018a) Evaluation of the seismic hazard and risk in the Batken region, Kyrgyzstan, with special attention to waste deposits. ICES technical report to MSF
- Wyss M, Gupta S, Rosset P (2018b) Casualty estimate in repeat Himalayan earthquakes in India. *Bull Seismol Soc Am* 108(5A):2877–2893
- Xu Y, Roecker SW, Wei R, Zhang W, Wei B (2006) Analysis of seismic activity in the crust from earthquake relocation in the central Tien Shan. *Bull Seismol Soc Am* 96:737–744
- Zhang P, Shen Z, Wang M, Gan W, Bürgmann R, Molnar P et al (2004) Continuous deformation of the Tibetan Plateau from global positioning system data. *Geology* 32:809–812



# Numerical modelling of the near-field velocity pulse-like ground motions of the Northridge earthquake

Quanbo Luo<sup>1</sup> · Feng Dai<sup>1</sup> · Yi Liu<sup>1</sup> · Mengtan Gao<sup>2</sup>

Received: 30 November 2019 / Accepted: 27 June 2020 / Published online: 6 July 2020  
© Institute of Geophysics, Polish Academy of Sciences & Polish Academy of Sciences 2020

## Abstract

The seismic records acquired during the 1994  $M_w$ 6.7 Northridge earthquake provide important data for studying the pulse-like ground motions in the vicinity of reverse faults. We selected 106 horizontal records from 468 strong ground motion records in the near-field region and rotated the original records into fault-parallel and fault-normal orientations. Large velocity pulses were simulated by the 3D finite difference method using a kinematic source model and a velocity structure model. Regression analysis was performed on the simulated and observed amplitudes of the velocity time history and response spectrum using the least-squares method. Our results show that the released energy and rupture time of asperities in the source model have important effects on the near-field velocity pulses, and the asperity near the initial rupture contributes more to the velocity pulses than does the asperity near the central region. The unidirectional and bidirectional characteristics of large velocity pulses are related to the thrust slip and rupture direction of the fault. The pulse period and the characteristic period are positively correlated with the rise time, and the pulse peak is regulated by multiple parameters of the subfaults. The distributions of the simulated PGV and Arias intensity agree well with the observed records, in which the contours exhibit asymmetric distribution and irregular elliptical attenuation in the near-field region, and the distributions exhibit a significant directivity along the fault. Moreover, the attenuation rate decreases with increasing distance from the fault. In addition, the fault-normal component is larger than that on the fault-parallel component, and the former decays faster. Velocity pulses larger than 30 cm/s are most likely to be distributed within approximately 15 km from the fault plane of the Northridge earthquake. Thus, the revealed pattern of the near-field velocity pulse-like ground motions indicates their close relation with the most severe earthquake effects.

**Keywords** Northridge earthquake · Finite difference method · Large velocity pulse · Source model · PGV · Response spectrum · Arias intensity

## Introduction

On 17 January 1994, at 4:31 local time (12:31 UTC), an earthquake of magnitude  $M_w$ 6.7 took place in the Northridge area northwest of Los Angeles, California. The epicentre was located in San Fernando Canyon at (34.206° N, 118.554° E) with a shallow focal depth of approximately 17.5 km, as shown in Fig. 1. This strong earthquake caused

many casualties and property losses, leading to more than 60 deaths and 9000 injuries, and a large number of high-rise buildings and bridges were damaged (Liu et al. 2012). The Northridge area is located on the West Coast of the USA within the largest seismic belt in the world, namely the circum-Pacific Ring of Fire, which displays a high incidence of earthquakes.

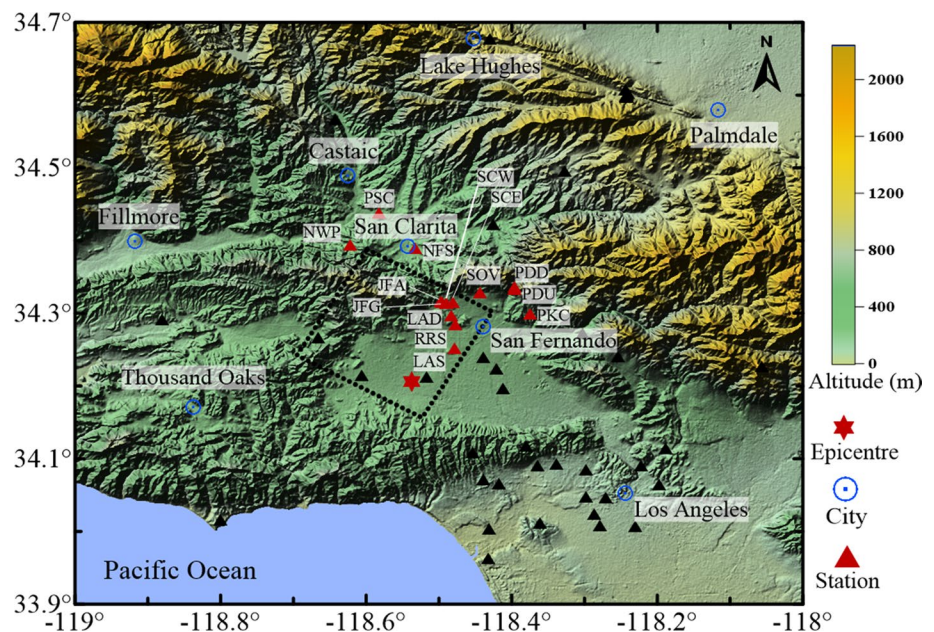
When a causative fault ruptures with a velocity close to the shear wave, the earthquake rapidly releases the enormous strain energy accumulated during the long-term tectonic movement. The large velocity pulses are characterized by high amplitude, long period, which arrive early at time histories as simple harmonic oscillations. Based on their characteristics, velocity pulses are divided into one-side and two-side pulses (Kawase and Aki 1990; Heaton et al. 1995; Oglesby and Archuleta 1997). These large velocity pulses

✉ Yi Liu  
liuyi\_scu@163.com

<sup>1</sup> State Key Laboratory of Hydraulics and Mountain River Engineering, College of Water Resource and Hydropower, Sichuan University, Chengdu, Sichuan 610065, China

<sup>2</sup> Institute of Geophysics, China Earthquake Administration, Beijing 100081, China

**Fig. 1** A topographic map of the Northridge area. The dashed rectangle depicts the surface projection of the causative fault plane for the Northridge earthquake. The strong ground motion stations with and without pulse records are indicated by triangles, and stations with and without pulse records are indicated by red triangles and black triangles, respectively. The surrounding cities of the Northridge earthquake are marked with the blue open circles. The epicentre is marked with a star



can cause substantial damage to large structures, as they can easily cause large inter-story displacements and permanent deformation (Bertero et al. 1978; Malhotra 1999; Li et al. 2020). In recent years, a small number of velocity pulses have been recorded during global earthquakes; for example, 29 pulses were recorded in the 1999  $M_w$ 7.6 Chi–Chi earthquake, 9 in the 2010  $M_w$ 7.0 Darfield earthquake, and 7 in the 2008  $M_w$ 7.9 Wenchuan earthquake. These earthquakes have attracted considerable interest in the fields of seismology and engineering. With the rapid development of the technology, buildings with higher natural vibration periods (large bridges, high-rise buildings, and oil storage tanks) are gradually proliferating. Therefore, the study of near-field long-period velocity pulses is of great significance for seismic hazard analysis and seismic design.

Because of the uncertainties in ground motions and the scarcity of seismographs, the Pacific Earthquake Engineering Research Center (PEER) has collected fewer than 200 pulse recordings, which is a rather poor sample to provide a statistical model of the characteristics of pulse-like ground motions. In order to compensate for the shortage of pulse records, models that can effectively simulate velocity pulses have been proposed by several researchers (Dickinson and Gavin 2011; Li 2016; Pu et al. 2017). However, models based on engineering approaches do not account for the rupture history. To cope with this, deterministic methods have been proposed to simulate the near-field velocity pulses emitted from large seismogenic sources.

For the simulation of time histories within the low period range of engineering interest ( $< 1$  s), the stochastic (Boore 2003; Motazedian and Atkinson 2005; Zhang and Yu 2010) and empirical Green's function (Irikura

1983; Choudhury et al. 2016) methods are usually used. Beresnev and Atkinson (1998a) performed a successful simulation of the acceleration histories that recorded the 1985  $M_w$  8.1 Mexico earthquake by using the stochastic finite-fault method. Li et al. (2017) simulated the acceleration records of the 1997 Kyushu earthquake by using the empirical Green's function method and analysed the relevant engineering parameters. The above methods are widely used for simulating short-period strong ground motions, but the simulation accuracy of near-field long-period ground motions is low (Irikura 1983; Li et al. 2018). Alternatively, for the low-frequency components (less than 1 Hz) in near-field ground motions, it is more suitable to apply a deterministic method.

Long-period ground motions can be effectively simulated by the 3D finite difference method (Kramer 1996; Graves 1998; Pitarka 1999; Luo et al. 2019). Many seismologists have verified the feasibility of the 3D finite difference method for simulating the near-field long-period ground motions of different earthquakes, the research results of which provide important guidance for disaster reduction. Gao et al. (2002) simulated the basin effect in Beijing and noted that the amplification factor of the local area is approximately 2. Maeda et al. (2016) performed a seismic hazard analysis of long-period ground motions generated by many scenarios of a megathrust earthquake in Nankai. Furumura et al. (2019) simulated the propagation of seismic waves in heterogeneous structures and forecasted the long-period ground motions generated by large earthquakes in sedimentary basins, and validated the effectiveness of the finite difference method by using observed waveform data from the 2007  $M_w$ 6.6 Niigata and 2011  $M_w$ 9.0 Tohoku earthquakes.

This study attempts to simulate the near-field large velocity pulses of the 1994 Northridge earthquake using the 3D finite difference method. This paper is organized as follows. Firstly, we describe the simulation method and the source function. Then, the near-field velocity pulses are identified from the seismic records of the Northridge earthquake. "Model and parameter setting" section establishes the kinematic source model and velocity structure model and presents the regional calculation parameters. In "Results and discussion" section, the characteristics and distribution fields of the near-field velocity pulse-like ground motions are illustrated through the comprehensive analysis and discussion of the numerical simulation results. The results can reveal the causes of large velocity pulses and help to analyse the responses of large-scale engineering structures to ground motions. "Conclusions" section summarizes the whole study.

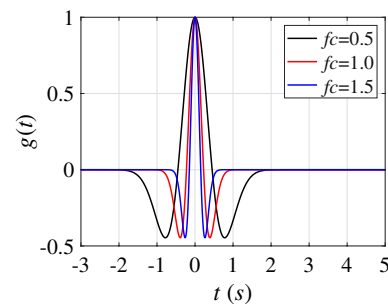
## Finite difference simulation method

Compared with the finite element method and the discrete wavenumber method, the finite difference method proposed by Aki (1968) can effectively simulate long-period ground motions in a large area consisting of an inhomogeneous medium. In the finite difference method, which has been continuously improved by seismologists over the decades (Mikumo et al. 1987; Aoi and Fujiwara 1999), the study area can be divided into discrete grids in the horizontal and vertical directions based on the different characteristics of geological structures, which greatly improves the calculation efficiency while ensuring the calculation accuracy. The relationship between the velocity pulses and the source model parameters in the near-field long-period ground motions can be studied utilizing the 3D finite difference method.

To simulate the long-period velocity pulses with the 3D finite difference method, it is necessary to establish a suitable source model including the geometric parameters and kinematic parameters. The fault plane is divided into finite discrete grids, and then the slip, seismic moment, and source time function are embedded into the velocity–stress difference equation to obtain the velocity history generated during the earthquake (Aoi et al. 2012; Maeda et al. 2014; Iwaki et al. 2016). The source function is the temporal and spatial function of each subfault during the rupture process. We use the Ricker wavelet to simulate the near-field long-period velocity pulses generated by the Northridge earthquake:

$$g(t) = (1 - 2\pi^2 f_c^2 t^2) \exp(-\pi^2 f_c^2 t^2) \quad (1)$$

where  $g(t)$  is the amplitude of the Ricker wavelet and  $f_c$  is the characteristic frequency, i.e. the reciprocal of the rise time of the subfault from initial rupture to slip termination. Figure 2 shows waveforms with characteristic frequencies



**Fig. 2** The waveforms corresponding to the source function with characteristic frequencies of 0.5, 1, and 1.5 Hz

of 0.5, 1, and 1.5 Hz. The wavelet function proposed by Ricker (1943) is widely used to simulate seismic waves (Ji et al. 2002; Wang 2015; Liu et al. 2016). The seismic wave received by a station on the surface is usually a short vibration that is excited by the subfault and propagates through the underlying medium.

## Strong motion recordings

The Next Generation Attenuation (NGA) database includes a large number of strong motion records from the Northridge earthquake, providing valuable fundamental data for studying near-field large velocity pulses. However, few stations are situated near the fault, and the spacing is variable; thus, a relatively small number of velocity pulses were recorded during this earthquake. Baker (2007) proposed three criteria for identifying velocity pulses: The pulse index in formula (2) is greater than 0.85, and the pulse appears early in the velocity time history and the peak ground velocity (PGV) is greater than 30 cm/s.

$$PI = 1/[1 + e^{-23.3+14.6(PGV_{ratio})+20.5(E_{ratio})}] > 0.85 \quad (2)$$

where PI is the pulse index,  $PGV_{ratio}$  is the ratio of the residual PGV to the original record after the velocity pulse is extracted,  $E_{ratio}$  is the ratio of the residual energy to the original record.

We selected 106 horizontal records in the study area from 468 strong motion records and identified 14 stations with velocity pulses and 39 stations without pulses based on the above criteria. For the stations that recorded velocity pulses during the Northridge earthquake, PKC and NWP are located at the ends of the fault, and 7 stations (LAS, RRS, LAD, SCE, SCW, JFA, and JFG) are located on the hanging wall, while the remaining stations are located on the footwall; the pulse data from these stations are listed in Table 1. Because the velocity waveforms recorded in the vertical direction do not meet the pulse standard and buildings

**Table 1** Basic information of the 14 stations that recorded velocity pulses

Abbrev.	Station Name	Lat. (°N)	Long. (°W)	PGV-FP (cm/s)	PGV-FN (cm/s)	ClstD (km)	Owner
PKC	Pacoima Kagel Canyon	34.296	118.375	29.5	56.3	7.26	CDMG
PDU	Pacoima Dam Upper Left	34.330	118.396	22.9	67.8	7.01	CDMG
PDD	Pacoima Dam Downstream	34.334	118.396	12.5	47.9	7.01	CDMG
LAS	LA-Sepulveda VA Hospital	34.249	118.479	48.1	47.1	8.44	USGS
SOV	Sylmar-Olive View Med FF	34.326	118.444	43.3	102.8	5.30	CDMG
RRS	Rinaldi Receiving Station	34.281	118.478	46.2	117.4	6.50	LADWP
LAD	LA Dam	34.294	118.483	51.5	74.7	5.92	LADWP
SCE	Sylmar-Converting Station East	34.312	118.481	66.3	90.5	5.19	LADWP
SCW	Sylmar-Converting Station West	34.311	118.490	81.5	121.0	5.35	LADWP
JFA	Jensen Filter Plant Administrative Bld.	34.312	118.496	90.9	104.3	5.43	USGS
JFG	Jensen Filter Plant Generator Bld.	34.313	118.498	58.2	69.3	5.43	USGS
NFS	Newhall-Fire Station	34.387	118.533	33.6	84.6	5.92	CDMG
PSC	Pardee-SCE	34.435	118.582	72.4	50.7	7.46	SCE
NWP	Newhall-W. Pico Canyon Rd.	34.391	118.622	69.0	108.4	5.48	USC

ClstD is the closest distance from the recording station to the ruptured area. Owner is the name of agency that collected the data

USGS U.S. Geological Survey, LADWP Los Angeles Department of Water and Power, CDMG California Division of Mines and Geology, USC University of Southern California, SCE Southern California Edison

are affected mainly by horizontal vibrations, this paper studies only the horizontal components of the near-field ground motions. Considering that the directivity effect and the source radiation have different influences on the horizontal components, to obtain a reference for a seismic comparison and to establish a relationship between the original seismic records and the strike of the fault, we rotate the original orthogonal horizontal components into fault-parallel (N122° E) and fault-normal (N212° E) orientations.

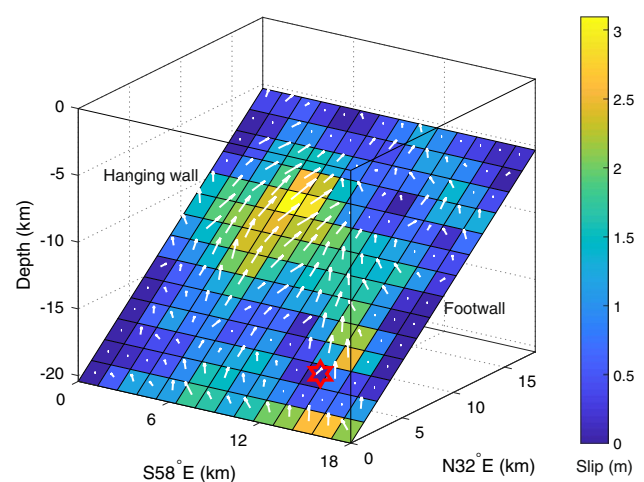
## Model and parameter setting

### Source model

The focal depths determined by the CMT Project, PEER, and USGS range from 16.8 to 18.2 km (with an average of approximately 17.5 km); the strike is S58° E, and the dip angle of the fault plane is approximately 40° to the southwest. To clarify the fault geometry and rupture motion characteristics of the Northridge earthquake, numerous scholars have performed considerable research. Zeng and Anderson (1996) obtained a composite source model of the earthquake using a genetic algorithm and indicated that a large amount of slip occurred near the central source region. Wald et al. (1996) combined teleseismic, strong motion, GPS displacement, and permanent uplift recordings to obtain the slip distribution characteristics of the Northridge earthquake. Beresnev and Atkinson (1998b) divided the fault plane into 20 subfaults and verified the slip distribution characteristics on the fault plane by using the simulated acceleration

histories at 28 rock sites. The above results show that the fault geometry can be represented by a rectangular plane with a complicated and inhomogeneous distribution of slip.

The Northridge earthquake was triggered by a blind causative fault (Wald et al. 1996; Ji et al. 2002). The rupture occurred along a thrust from approximately 20 km to 5 km below the surface at a dip angle of 40° and was truncated by the San Fernando fault (Mori et al. 1995). We have established a source model formed by a rectangular plane in this study (Fig. 3). The fault ruptured 18 km along the strike approximately 5 km below the surface and ruptured



**Fig. 3** The fault model of the Northridge earthquake. The fault plane is divided into 196 subfaults, each showing the direction of the average slip with an arrow and the magnitude of the average slip with colour. The hypocentre is indicated by a star



downward approximately 24 km along dip. The projection of the fault plane on the surface forms the black dashed rectangle shown in Fig. 1. We divided the entire fault plane into 196 subfaults of  $1.286 \times 1.714$  km.

The energy generated by the rupture of an asperity during an earthquake greatly contributes to the strong ground motion; thus, the source model of an asperity is significant for evaluating the seismic effect on an engineering structure (Kamae and Irikura 1998). The strength of the asperity region is less than the stress field, thereby enhancing the fault rupture, which experiences a high stress drop during the rupture of the fault (Aki 1984). We extracted the position, quantity, and area of the asperities from the inhomogeneous slip distribution and assumed two asperities for the Northridge earthquake (Fig. 4): small asperity A is approximately  $19.8 \text{ km}^2$ , and large asperity B is approximately  $72.7 \text{ km}^2$ . Somerville et al. (1999) studied the spatial slip distribution of 15 crustal earthquakes with magnitudes greater than 5.7 worldwide and proposed that the area ratio of asperities to the entire fault is 0.22, and Murotani et al. (2008) proposed that the area ratio of a plate boundary earthquake is close to 0.2. In this study, the ratio of the total area of both asperities to the area of the entire fault is approximately 0.21, which is basically consistent with Somerville et al. (1999). In previous studies, i.e. the 1997  $M_w$ 6.0 Kagoshima, 2000  $M_w$ 6.6 Tottori, and 2004  $M_w$ 6.6 Chuetsu earthquakes (Irikura and Miyake 2011; Iwaki et al. 2016), asperities have been

approximated by a rectangle. However, the large slip on the fault plane is not necessarily located in a rectangular area. Based on the spatial inhomogeneity of the slip distribution (Wald et al. 1996), we set asperities A and B of the Northridge earthquake to have rectangular and irregular shapes, respectively.

The seismic moment is used to measure the energy released by an earthquake and thus has an important influence on the ground motion. The distribution of the seismic moment on the fault plane is shown in Fig. 4. This study determined that the total seismic moment of the Northridge earthquake is  $1.15 \times 10^{19} \text{ N m}$ , which is close to  $1.3 \pm 0.2 \times 10^{19} \text{ N m}$  estimated by Wald et al. (1996). The seismic moment of the asperities and background region are distributed according to formula (3) proposed by Somerville et al. (1999), where the total seismic moment of both asperities is approximately  $7.47 \times 10^{18} \text{ N m}$ , and that of the background region is  $4.03 \times 10^{18} \text{ N m}$ .

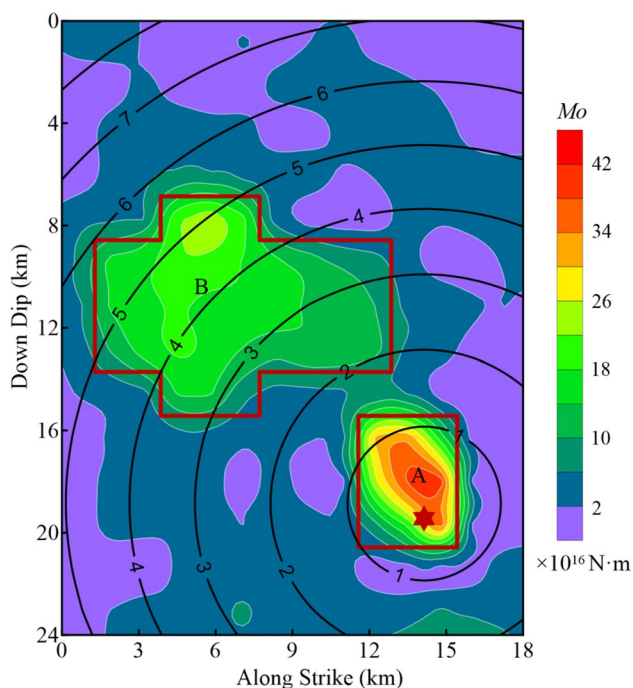
$$M_{oa} = \mu D_a S_a \quad (3)$$

where  $\mu$  is the average crustal rigidity, its value is about 30 Gpa.  $M_{oa}$  is the seismic moment of the asperity, and  $S_a$  is the area of the asperity.  $D_a$  is the average slip of the asperity, and its value (approximately 2.7 m) is the total slip on the asperity divided by the number of subfaults.

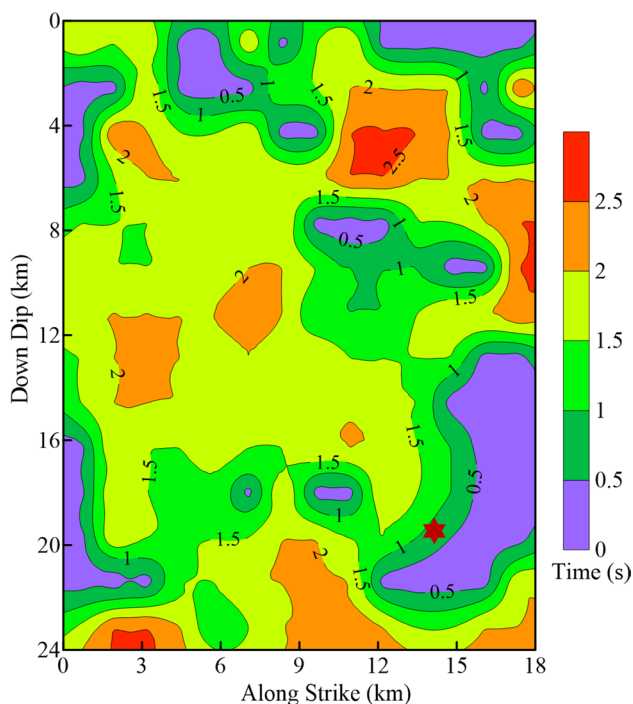
In our model, the Northridge earthquake began with a circular rupture near the bottom of the fault plane that propagated from the southeast to the northwest with an average rupture velocity of 2.8 km/s. Field et al. (1998) studied the nonlinear sediment response during the Northridge earthquake using a uniform circular rupture pattern, and the results indicated the effectiveness of circular rupture. Hartzell et al. (1996) found that the fault rupture velocity at the early stage of the earthquake was 2.8–3.0 km/s, while the velocity after 3 s was 2.0–2.5 km/s. We used a varying rupture velocity for the source rupture pattern: the velocity rapidly decreased outward from 3.0 km/s in the nucleation zone to 2.5 km/s over a total rupture time of approximately 8 s. The rise time of the fault slip is inhomogeneously distributed (Hartzell et al. 1996; Wald et al. 1996), and the nucleation zone is relatively small at approximately 0.6 s, although the rise time tends to increase outward, as shown in Fig. 5.

### Velocity structure model

The crustal velocity structure reflects the stratigraphic sequence from the surface to the Moho and the variation in the seismic velocity with depth. A reasonable velocity structure model, which has an important influence on the simulation results of long-period velocity pulses, can be established according to the changes in the physical properties of each layer. In the numerical simulation of near-field



**Fig. 4** The seismic moment distribution of the Northridge earthquake. The asperities are surrounded by red lines. The black lines are contours of rupture time at 1 s intervals



**Fig. 5** The rise time distribution of the fault displacement; the adjacent contours are separated by 0.5 s intervals

velocity pulses, the viscoelastic effect of the crustal medium needs to be considered. When the S-wave velocity is less than 1–2 km/s, the ratio of the attenuation factor  $Q$  to  $V_s$  is close to 0.02, and when the S-wave velocity is greater than 2 km/s, the ratio in the Los Angeles area is approximately 0.1 (Olsen et al. 2003). We set the viscoelastic properties of the velocity model according to existing studies in the region (Magistrale et al. 1992; Olson et al. 1984, 2003). The adopted velocity model is presented in Table 2.

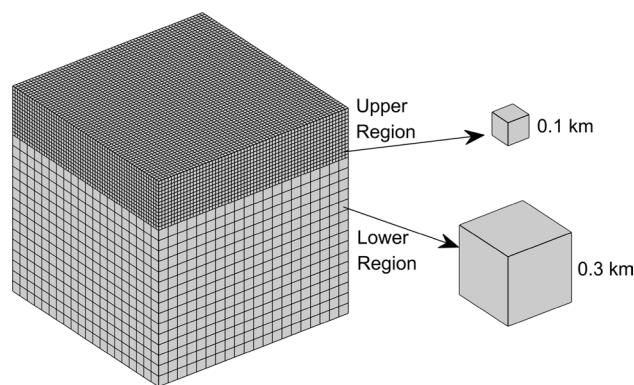
The crust in the Northridge area is divided into grids with an interface of 8 km beneath the surface because the shallow strata in the region have smaller seismic velocities than the deep strata. The entire region is divided into small grids, considerable amounts of computational time and memory will be consumed. However, the region is divided into large grids, and the simulation results may numerically diverge. To satisfy the accuracy and efficiency of calculation at the

same time, we used a non-uniform grid to divide the study region (Fig. 6). In the upper and lower regions, the grid spacing is 0.1 km and 0.3 km respectively, with a total of approximately  $6.84 \times 10^7$  grids. To ensure the stability of the numerical simulation, five grids were used in one wavelength under the condition of a fourth-order precision. At the same time, the simulated low frequency was appropriately extended to a high frequency, and the upper limit of the frequency for the velocity pulses simulation was taken as 1.4 Hz. The detailed calculation parameters are listed in Table 3.

## Results and discussion

### Waveform comparison

The simulated and observed waveforms of the 28 velocity pulse histories of the near-field ground motions are shown in Fig. 7. The red dashed lines indicate the simulated results, the black solid lines indicate the observed records, and all the data are low-pass filtered with a cut-off frequency of 1.4 Hz. Most velocity histories match well regarding the amplitude and phase, and fewer pulses are recorded on the fault-parallel (FP) component (displaying complex



**Fig. 6** Schematic diagram of 3D non-uniform grid configuration for the local region

**Table 2** Relevant velocity structure model for the Northridge area

Depth (km)	Thickness (km)	$V_p$ (km/s)	$V_s$ (km/s)	Density ( $\text{kg/m}^3$ )	$Q$
0	0.5	2.1	1.08	2100	22
0.5	3.5	4.0	2.15	2500	215
4.0	2.5	4.8	2.65	2600	265
6.5	14.0	6.1	3.50	2900	350
20.5	14.5	7.0	4.00	3000	400
35.0	$\infty$	7.8	4.50	3300	450

**Table 3** Calculation parameters used in this study

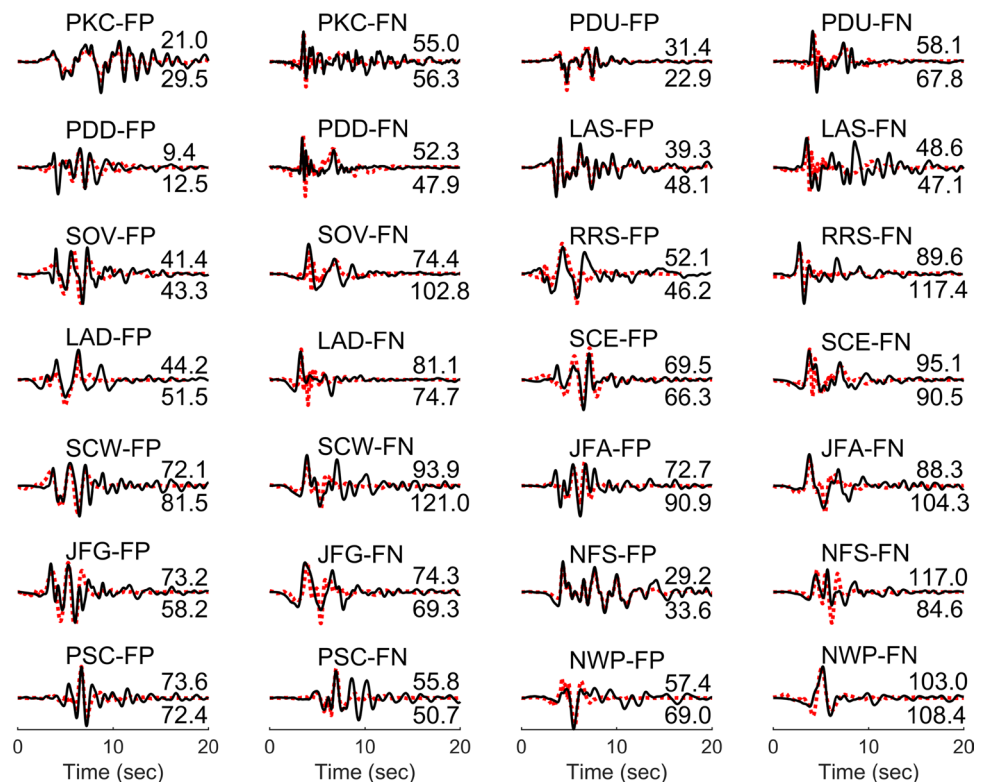
Model size (nx × ny × nz)	88.2 × 92.4 × 35
Total time steps	4000
Time step (s)	0.005
Upper limit frequency (Hz)	1.4
Total grid points	6.84 × 10 <sup>7</sup>
Receiving stations	53
Simulation area range	33.9° N–34.7° N, 118° W–119° W

waveforms) than on the fault-normal (FN) component (displaying simple waveforms). At the rupture front-end station NWP, the velocity waveforms exhibit one-side long-period velocity pulses, the pulse peak of the fault-normal component is greater than 100 cm/s, and the pulse period is approximately 2 s, while at the rupture back-end station PKC, the pulse peak and period of are approximately half of those at station NWP, and the fault-parallel component did not record the velocity pulse history. The velocity pulse is affected by the seismic Doppler effect and the orientation of the station, the energy of the rupture radiation of each sub-fault is stacked at the front end of the rupture, while the time it takes for the energy to reach the back end of the rupture is delayed; thus, the velocity histories of stations NWP and PKC are characterized by forward directivity and backward directivity effects, respectively.

Among the 7 stations on the hanging wall, station LAS is the farthest (8.44 km) from the fault plane; at this station, the pulse peaks of the two components are the smallest and basically equivalent, and the fault-parallel component has a more obvious two-side pulse than the fault-normal component in the velocity histories, which indicates that the fling-step effect has a greater influence on the vicinity of station LAS than does the directivity effect. Mavroeidis and Papa-georgiou (2003) argued that the peak of the near-field pulse does not increase indefinitely, but there is a typical threshold of approximately 100 cm/s. In the actual recordings, the pulse peaks of the fault-normal component for stations RRS and SCW are approximately 120 cm/s. Compared with the fault-parallel components, the fault-normal components for stations RRS and LAD have larger peaks and smaller pulse periods, and the velocity pulses are also more significant. These differences are related to the positions of the stations on the active hanging wall and are greatly affected by the fling-step effect caused by the thrust motion along the fault. Stations SCE, SCW, JFA, and JFG are located at similar positions, but the actual recorded pulses are different, which is related to the local site of each station. For example, stations SCE and SCW are located on rock and soil, respectively; thus, the velocity history of SCW records a larger pulse peak than that of SCE, reflecting the amplification effect of the soil layer on the pulse peak.

On the footwall, the three stations (PDU, PDD, and SOV) near the initial rupture end of the fault are closer to asperity

**Fig. 7** Comparison of the simulated (red dashed lines) and observed (black solid lines) waveform for the 28 velocity pulses. The station abbreviation along with the component name is shown above each curve. The maximum amplitudes in cm/s are shown to the right of the curves, simulated value is indicated above the end of each curve, and observed value is indicated below the end of each curve. The strong ground motion data applied low-pass filtering at 1.4 Hz



A than the two stations (NFS and PSC) near the front end of the rupture, and the former three stations exhibit more pronounced velocity pulses than the latter two, while the high amplitudes after the pulses on their velocity histories are mainly affected by the asperity B. Station SOV is the closest (5.3 km) to the fault plane with a recorded pulse peak greater than 100 cm/s on the fault-normal component, and the waveform has a large wave period after the initial pulse, which is related to the inhomogeneous distribution of the rise time on the fault plane. The small rise times of the nucleation zone produce short-period and high-amplitude velocity pulses, while the large rise times on the fault plane produce long-period and low-amplitude waveforms.

In the numerical simulation of velocity pulses, we found that the transverse component of the S-wave is larger than the radial component of the P wave due to the radiation of the ruptures on the subfaults, and thus, the fault-normal component has a larger pulse than the fault-parallel component at most near-field stations. Comparing the pulse peaks between the hanging wall and footwall stations, there are three velocity pulses (RRS-FN, SCW-FN, and JFA-FN) that exceed 100 cm/s on the hanging wall, while only SOV-FN recorded a velocity pulse of approximately 100 cm/s on the footwall. At the same time, station SCW on the hanging wall displays a larger pulse peak than station SOV on the footwall at a similar distance from the fault. Therefore, the stations on the active hanging wall are more affected by the asperities on the fault plane and more easily record the velocity pulses. It can be seen from the simulation results of all the stations that recorded the velocity pulses of the Northridge earthquake that the 3D finite difference method can effectively simulate long-period velocity pulses, but some of the short-period velocity waveforms are not ideal. A large amount of data was recorded throughout the near-field region of the Northridge earthquake, but the distribution of strong motion stations with pulse records was unbalanced, and the number of stations near the front end of the rupture was much smaller than that near the back end. It is therefore necessary to analyse the characteristics of the velocity pulses from the spatial distribution of near-field ground motions.

### PGV analysis

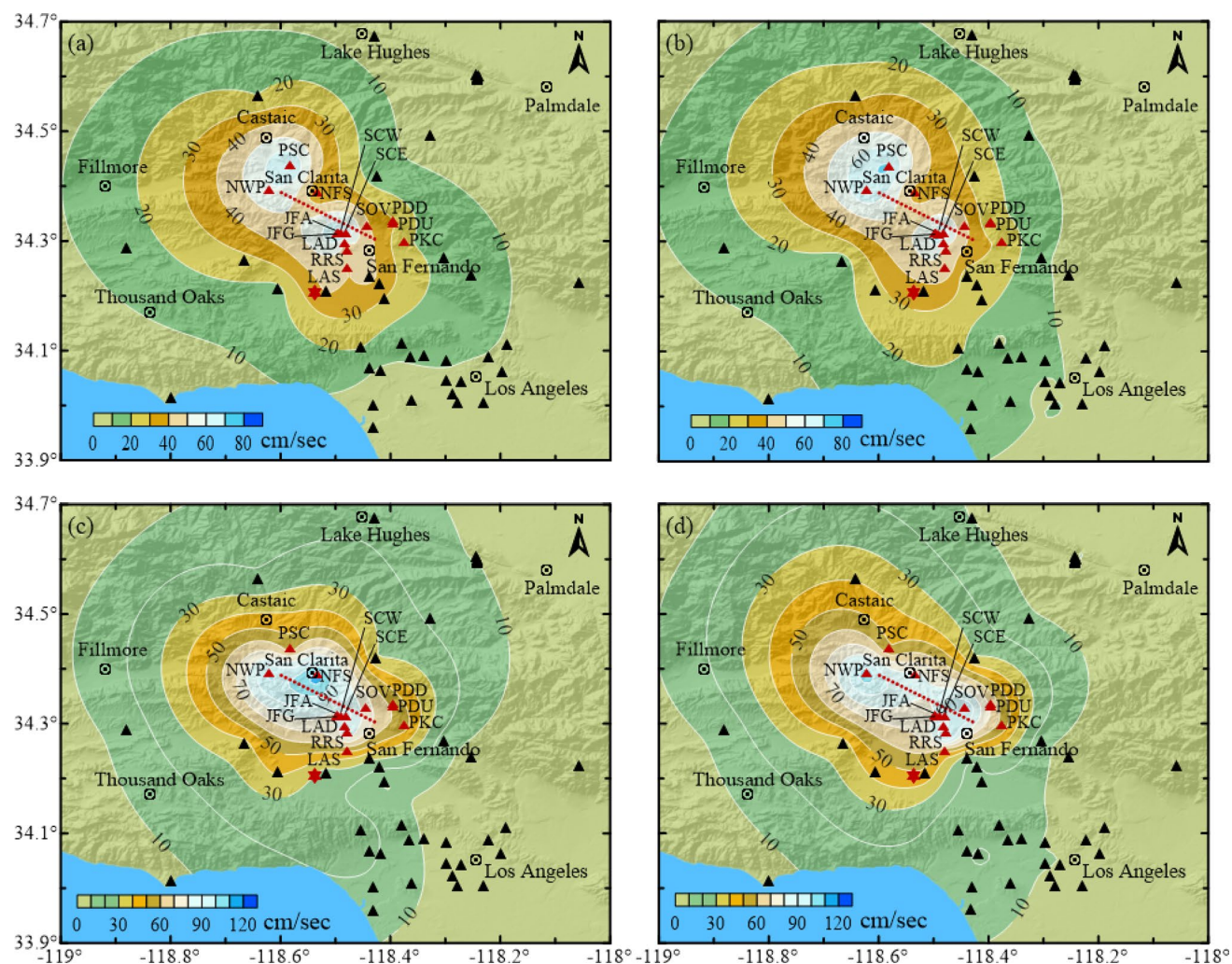
The peak ground velocity (PGV) is one of the most important parameters reflecting the intensity of ground motion. It can provide a good reference for estimating the seismic intensity, determining seismic zoning, and future urban planning. The damage attributable to near-field strong motion is mainly related to long-period components, and the peak velocity is more likely than the peak acceleration to reflect the long-period characteristics of near-field ground motions (Wald et al. 1999; Xu and Xie 2005).

To analyse the distribution characteristics of the long-period PGV in the near-field region, we plot contour maps with an interval of 10 cm/s by using the simulated peaks from 53 stations (no pulses were recorded at 39 stations) and perform low-pass filtering with a cut-off frequency of 1.4 Hz for all the data; the simulated long-period PGV distribution is then compared with the actually observed records, as shown in Fig. 8. The simulated PGV is similar to the observed PGV with distribution characteristics along the fault strike. The PGV at the front end of the fault rupture has a wider distribution than that at the back end of the rupture, the former exhibits slower decay, and the near-field ground motion reflects the typical directivity effect. The peak ground velocities are similar between the horizontal components, but the intensity and attenuation of each component are different. The PGV on the fault-parallel component is significantly smaller than that on the fault-normal component, and velocity pulses are recorded more frequently on the fault-normal component. It can be seen from the spacing and intensity of the contours that the PGV decays faster in the vicinity of the fault, the decay rate of the PGV gradually decreases as the fault distance increases, and the fault-parallel component decays slower than the fault-normal component.

Strong ground motions are mainly concentrated in the vicinity of stations SCW and NWP, and the maximum peaks on the fault-parallel and fault-normal components are approximately 90 cm/s and 120 cm/s, respectively. The simulated value of the near-field long-period PGV is slightly smaller than the observed value in the local area. The reasons for this difference may include the uncertainties in the source parameters, the seismic wave disturbances caused by changes in the terrain, and the amplification of the ground motions in the Los Angeles Basin and San Fernando Basin.

### Pseudo-velocity response spectra comparison

The pseudo-velocity response spectrum is presented as the maximum pseudo-velocity response curve of a single-degree-of-freedom elastic system that changes with the natural vibration period under a given ground motion. The response spectrum derives from the combination of structural dynamic characteristics (the natural vibration period, vibration mode, and damping) and ground motion; accordingly, the resonance effect of a structure in an earthquake can be calculated by the response spectrum. The characteristic period discussed in this paper refers to the period corresponding to the maximum amplitude of the pseudo-velocity response spectrum; the near-field velocity pulse of the Northridge earthquake has a large characteristic period, resulting in serious damage to long-period large-scale structures, especially lifeline engineering and building structures in the near-field region. Therefore, the characteristic period



**Fig. 8** Contour maps of PGV in cm/s obtained from 53 stations surrounding the fault. All the data are low-pass filtered by a frequency of 1.4 Hz. **a** Simulated values of the fault-parallel components; **b**

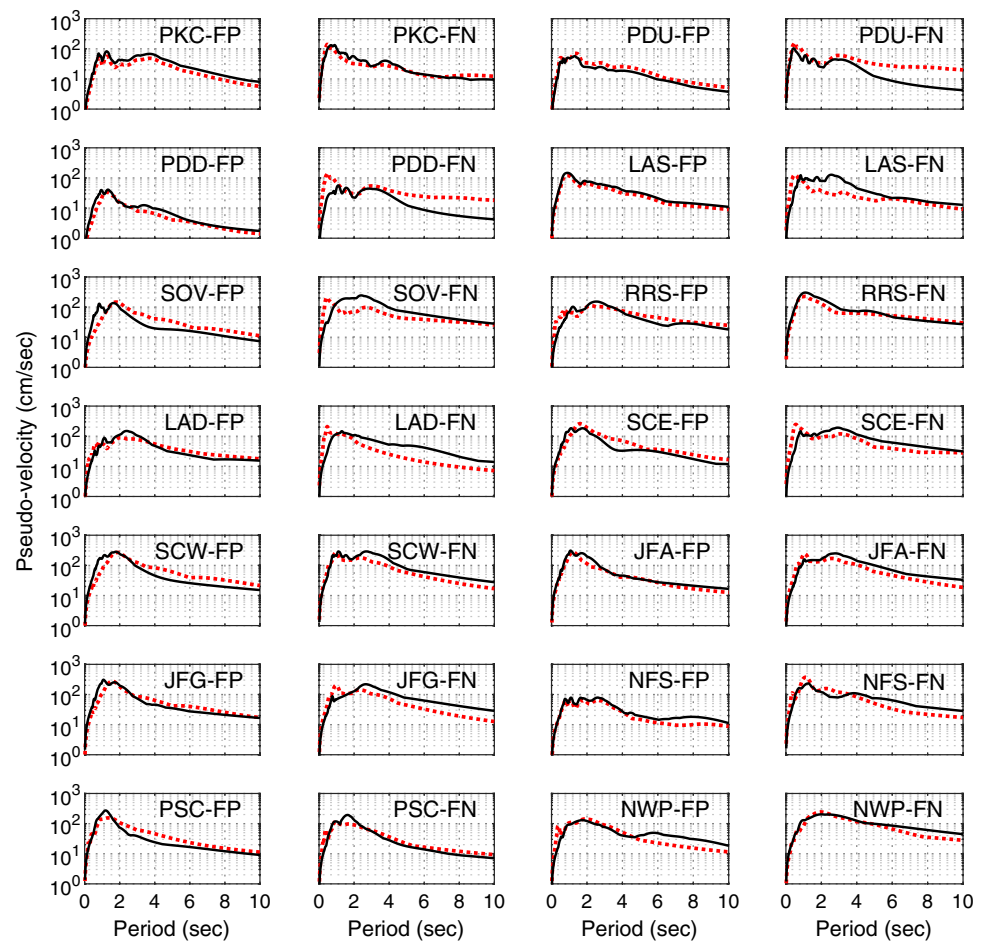
observed values of the fault-parallel components; **c** simulated values of the fault-normal components; **d** Observed values of the fault-normal components

of the pseudo-velocity response spectrum is of great significance for engineering research.

The pseudo-velocity response spectra for the horizontal components of each station that recorded velocity pulses during the Northridge earthquake are shown in Fig. 9. The damping ratio is 5%, the period is 1–10 s, the red dashed line indicates the simulated response spectrum, and the black solid line indicates the observed response spectrum. There are some differences among the characteristic periods of the pseudo-velocity response spectra. For example, the characteristic period of the response spectrum of stations SOV, SCW, and NWP is approximately 2 s, and the characteristic period of stations PDU, LAS, and PSC is approximately 1 s, while the characteristic periods of stations RRS, LAD, and JFA on different components differ by approximately 1 s. Comparing the pseudo-velocity response spectra with the velocity histories, it can be determined that the characteristic

period of the response spectrum is positively correlated with the pulse period of the velocity history. The maximum spectral value of the pseudo-velocity response spectrum also shows some differences between the different components of each station. For example, the maximum spectral values of station NWP on the fault-parallel and fault-normal components are approximately 130 cm/s and 200 cm/s, respectively. Therefore, the maximum spectral value of the pseudo-velocity response spectrum is related to the pulse peak of the velocity history. From the overall comparison of the pseudo-velocity response spectra, it can be seen that the simulated values are close to the observed values; however, the simulated and observed spectra of stations PDU and PDD on the fault-normal component display some differences after the characteristic period, which is more likely to be affected by local site effects compared with the long-period surface waves.

**Fig. 9** Comparison of the simulated (red dashed lines) and observed (black solid lines) pseudo-velocity response spectrum for the 28 pulses. The damping value is 5%

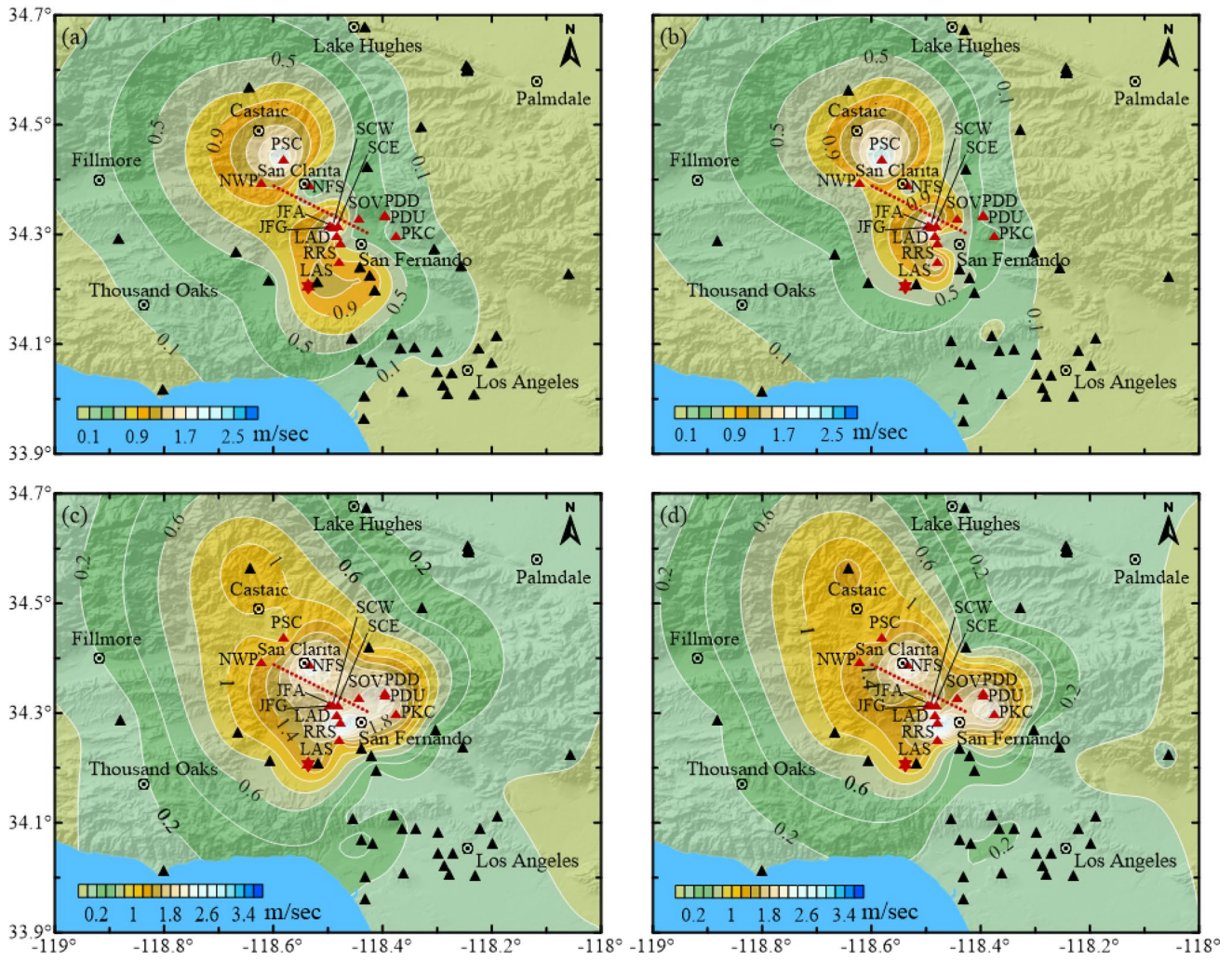


## Arias intensity analysis

The distribution characteristics of the Arias intensity in the near-field region of the Northridge earthquake are shown in Fig. 10. The intensity decays with increasing fault distance, and the intensity in the vicinity of the fault decays faster than that in the region far from the fault; moreover, the area of the Arias intensity at the front end of the fault rupture is wider than that at the back end, and the Arias intensity on the fault-parallel component is smaller than that on the fault-normal component. There are some differences in the concentrated areas of large intensity values around the fault. Large values are concentrated around stations PSC and SCW on the fault-parallel component and around stations NFS, RRS, and PDU on the fault-normal component; these large value areas are also indicative of large earthquake disasters. The simulated values of the Arias intensity are basically consistent with the observed values, but there are differences in some local areas far from the fault, which may be related to the whole earthquake procedure, as well as changes in the topography and local site conditions.

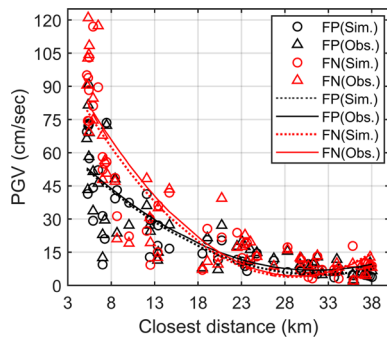
## Regression analysis

To verify the numerical simulation results of the near-field velocity pulses by the 3D finite difference method, regression analysis is performed on the simulated and observed values based on the least-squares method. The PGV and regression results for the horizontal components of 53 stations in the near-field region are shown in Fig. 11. The simulated values are similar to the observed values, and the PGV of each component exhibits a different attenuation trend with decreasing distance to the fault. The PGV on the fault-parallel component is smaller than that on the fault-normal component. The attenuation of the fault-parallel component occurs more slowly (i.e. the absolute slope of the regression line is small) within 18 km from the fault, but the attenuation of the two components occurs at a similar rate at distances greater than 18 km from the fault, which is basically consistent with the distribution characteristics of the PGV contours. One of the criteria for a velocity pulse that must be satisfied is a PGV greater than 30 cm/s. The velocity pulses on the fault-parallel and fault-normal components are most likely to be within 13 km and 15 km, respectively, of the fault, which



**Fig. 10** Arias intensity in m/s distribution obtained from 53 stations surrounding the fault. **a** Simulated values of the fault-parallel components; **b** observed values of the fault-parallel components; **c** simu-

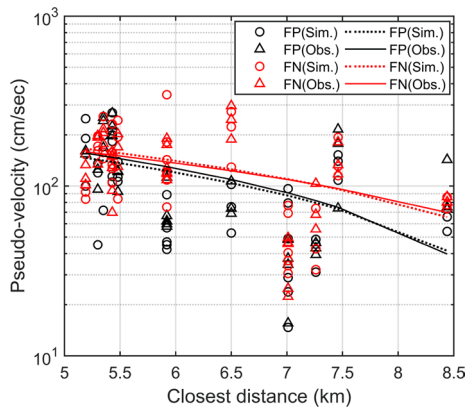
lated values of the fault-normal components; **d** observed values of the fault-normal components



**Fig. 11** Variations of PGV with the closest distance to fault plane for 53 strong ground stations. The simulated and observed values of the fault-parallel components are represented by black open circles and black open triangles, respectively. The simulated and observed values of the fault-normal components are represented by red open circles and red open triangles, respectively. The comparison of the regressions is represented by lines, where the simulated and observed values are indicated by the dashed and solid lines, respectively

also suggests that the source radiation effect produces large amplitudes along the direction perpendicular to the fault.

The characteristic period of the pseudo-velocity response spectrum is basically in the range of 1–2 s, and the spectral values corresponding to the vibration periods of 1, 1.5, and 2 s are taken from the response spectra. The simulated and observed values are also regressed, as shown in Fig. 12. It can be seen that the simulated results agree well with the observed records; the spectral values on the fault-normal component are larger than those on the fault-parallel component, and the attenuation of the former occurs more slowly than that of the latter, so the near-field velocity pulses are most likely to cause damage to structures along the direction perpendicular to the fault. In the seismic design of building structures, not only the short-period ground motions excited by active faults but also the long-period pulse-like ground motions should be considered. Clearly, the study of



**Fig. 12** Comparison of pseudo-velocity regression with the closest distance to fault plane for the 28 pulses

near-field velocity pulses is of great significance for earthquake prevention and disaster reduction.

## Conclusions

The near-field long-period velocity pulses of the Northridge earthquake were simulated by the 3D finite difference method. The simulated velocity histories, PGV distribution, pseudo-velocity response spectra, and Arias intensities were compared with the real observed values, and regression analysis verified the feasibility of simulating the near-field velocity pulses by this method. The simulation results are expected to apply to near-field pulse-like ground motion assessments, seismic hazard analysis, and the study of non-linear structural responses. The following conclusions can be drawn:

1. The source model affects the characteristics of the near-field long-period velocity pulses and the distribution of pulse-like ground motions. The rectangular asperity provides an important contribution to the pulse peaks, and the irregular asperity mainly affects the waveforms after the velocity pulses. One-side pulses on the fault-normal component are mainly affected by thrust slip and two-side pulses on the fault-parallel component affected by rupture direction. The velocity pulses on the fault-normal component are more abundant than those on the fault-parallel component; besides, the pulse period is positively correlated with the rise time, and the pulse peak is regulated by the seismic moment, the amount of slip, and the rise times on the subfaults. Some peaks exceed 100 cm/s.
2. The PGV contours exhibit an asymmetrical distribution in the near-field region, and the distribution at the front end of the rupture is larger than that at the back end.

The PGV exhibits irregular elliptical attenuation, and the PGV decay rate gradually decreases with increasing distance from the fault. Similar to the PGV distribution, the Arias intensity also exhibits a significant directivity effect and attenuation trend; the fault-normal component is larger than that on the fault-parallel component; and the intensity on the former decays faster than that on the latter, while the distributions of the maximum values on different components do not necessarily coincide.

3. The characteristic period of the pseudo-velocity response spectrum is basically in the range of 1–2 s, and the characteristic period is related to the pulse period. Velocity pulses greater than 30 cm/s are most likely to be distributed within approximately 15 km of the fault; in addition, the fault-normal component has a larger distribution range than the fault-parallel component. Hence, the near-field region should be considered an important area during the seismic design of building structures.

**Acknowledgements** The authors thank the financial support from the National Natural Science Foundation of China (No. 51779164) and the Youth Science and Technology Innovation Research Team Fund of Sichuan Province (2020JDTD0001). The strong motion records come from the Pacific Earthquake Engineering Research (PEER) Center NGA-West2 database. We appreciated the anonymous reviewer for their constructive comments and suggestions that greatly improve this manuscript.

## References

- Aki K (1968) Seismic displacements near a fault. *J Geophys Res* 73(16):5359–5376
- Aki K (1984) Asperities, barriers, characteristic earthquakes and strong motion prediction. *J Geophys Res Solid Earth* 89(B7):5867–5872
- Aoi S, Fujiwara H (1999) 3D finite difference method using discontinuous grids. *Bull Seismol Soc Am* 89(4):918–930
- Aoi S, Maeda T, Nishizawa N, Aoki T (2012) Large-scale ground motion simulation using GPGPU. AGU, fall meeting, S53G-06
- Baker JW (2007) Quantitative classification of near-fault ground motions using wavelet analysis. *Bull Seismol Soc Am* 97(5):1486–1501
- Beresnev IA, Atkinson GM (1998a) FINSIM—a FORTRAN program for simulating stochastic acceleration time histories from finite-faults. *Seismol Res Lett* 69(1):27–32
- Beresnev IA, Atkinson GM (1998b) Stochastic finite-fault modeling of ground motions from the 1994 Northridge, California, earthquake. I. Validation on rock sites. *Bull Seismol Soc Am* 88(6):1392–1401
- Bertero VV, Mahin SA, Herrera RA (1978) Aseismic design implications of near-fault San Fernando earthquake records. *Earthq Eng Struct D* 6(1):31–42
- Boore DM (2003) Simulation of ground motion using the stochastic method. *Pure Appl Geophys* 160(3–4):635–676
- Choudhury P, Chopra S, Roy KS, Sharma J (2016) Ground motion modelling in the Gujarat region of Western India using empirical Green's function approach. *Tectonophysics* 675:7–22



- Dickinson BW, Gavin HP (2011) Parametric statistical generalization of uniform-hazard earthquake ground motions. *J Struct Eng* 137(3):410–422
- Field EH, Zeng Y, Johnson PA, Beresnev IA (1998) Nonlinear sediment response during the 1994 Northridge earthquake: observations and finite source simulations. *J Geophys Res Solid Earth* 103(B11):26869–26883
- Furumura T, Maeda T, Oba A (2019) Early forecast of long-period ground motions via data assimilation of observed ground motions and wave propagation simulations. *Geophys Res Lett* 46(1):138–147
- Gao MT, Yu YX, Zhang XM, Wu J, Hu P, Ding YH (2002) Three-dimensional finite-difference simulations of ground motions in the Beijing area. *Earthq Res Chin* 18(4):356–364 (in Chinese)
- Graves RW (1998) Three-dimensional finite-difference modeling of the San Andreas fault: source parameterization and ground-motion levels. *Bull Seismol Soc Am* 88(4):881–897
- Hartzell S, Liu P, Mendoza C (1996) The 1994 Northridge, California, earthquake: investigation of rupture velocity, risetime, and high-frequency radiation. *J Geophys Res Solid Earth* 101(B9):20091–20108
- Heaton TH, Hall JF, Wald DJ, Halling MW (1995) Response of high-rise and base-isolated buildings to a hypothetical  $M_w$  7.0 blind thrust earthquake. *Science* 267(5195):206–211
- Irikura K (1983) Semi-empirical estimation of strong ground motions during large earthquakes. *Bull Disaster Prev Res Inst Kyoto Univ Jpn* 33(2):63–104
- Irikura K, Miyake H (2011) Recipe for predicting strong ground motion from crustal earthquake scenarios. *Pure Appl Geophys* 168(1–2):85–104
- Iwaki A, Maeda T, Morikawa N, Miyake H, Fujiwara H (2016) Validation of the recipe for broadband ground-motion simulations of Japanese crustal earthquakes. *Bull Seismol Soc Am* 106(5):2214–2232
- Ji C, Wald DJ, Helmberger DV (2002) Source description of the 1999 Hector Mine, California, earthquake, part I: wavelet domain inversion theory and resolution analysis. *Bull Seismol Soc Am* 92(4):1192–1207
- Kamae K, Irikura K (1998) Source model of the 1995 Hyogo-ken Nanbu earthquake and simulation of near-source ground motion. *Bull Seismol Soc Am* 88(2):400–412
- Kawase H, Aki K (1990) Topography effect at the critical SV-wave incidence: possible explanation of damage pattern by the Whittier Narrows, California, earthquake of 1 October 1987. *Bull Seismol Soc Am* 80(1):1–22
- Kramer S (1996) *Geotechnical earthquake engineering*. Prentice Hall, Upper Saddle River, NJ, pp 50–300
- Li XX (2016) Study on extraction of the velocity pulse and effects of inclusion on ground motion. Institute of Engineering Mechanics, China Earthquake Administration, Beijing, pp 10–11 (in Chinese)
- Li Z, Chen X, Gao M, Jiang H, Li T (2017) Simulating and analyzing engineering parameters of Kyushu earthquake, Japan, 1997, by empirical Green function method. *J Seismol* 21(2):367–384
- Li Z, Gao M, Jiang H, Chen X, Li T, Zhao X (2018) Sensitivity analysis study of the source parameter uncertainty factors for predicting near-field strong ground motion. *Acta Geophys* 66(4):523–540
- Liu T, Luan Y, Zhong W (2012) A numerical approach for modeling near-fault ground motion and its application in the 1994 Northridge earthquake. *Soil Dyn Earthq Eng* 34(1):52–61
- Liu J, Zhang J, Sun Y, Zhao T (2016) Comparison of methods for seismic wavelet estimation. *Prog Geophys* 31(2):0723–0731 (in Chinese)
- Li A, Liu Y, Dai F, Liu K, Wei M (2020) Continuum analysis of the structurally controlled displacements for large-scale underground caverns in bedded rock masses. *Tunn Undergr Space Technol* 97:103288
- Luo Q, Chen X, Gao M, Li Z, Zhang Z, Zhou D (2019) Simulating the near-fault large velocity pulses of the Chi-Chi ( $M_w$  7.6) earthquake with kinematic model. *J Seismol* 23(1):25–38
- Maeda T, Morikawa N, Iwaki A, Aoi S, Fujiwara H (2014) Simulation-based hazard assessment for long-period ground motions of the Nankai Trough megathrust earthquake. AGU, fall meeting, S31C-4438
- Maeda T, Iwaki A, Morikawa N, Aoi S, Fujiwara H (2016) Seismic-hazard analysis of long-period ground motion of megathrust earthquakes in the Nankai trough based on 3D finite-difference simulation. *Seismol Res Lett* 87(6):1265–1273
- Magistrale H, Kanamori H, Jones C (1992) Forward and inverse three-dimensional P wave velocity models of the southern California crust. *J Geophys Res Solid Earth* 97(B10):14115–14135
- Malhotra PK (1999) Response of buildings to near-field pulse-like ground motions. *Earthq Eng Struct D* 28(11):1309–1326
- Mavroeidis GP, Papageorgiou AS (2003) A mathematical representation of near-fault ground motions. *Bull Seismol Soc Am* 93(3):1099–1131
- Mikumo T, Hirahara K, Miyatake T (1987) Dynamical fault rupture processes in heterogeneous media. *Tectonophysics* 144(1):19–36
- Mori J, Wald DJ, Wesson RL (1995) Overlapping fault planes of the 1971 San Fernando and 1994 Northridge, California earthquakes. *Geophys Res Lett* 22(9):1033–1036
- Motazedian D, Atkinson GM (2005) Stochastic finite-fault modeling based on a dynamic corner frequency. *Bull Seismol Soc Am* 95(3):995–1010
- Murotani S, Miyake H, Koketsu K (2008) Scaling of characterized slip models for plate-boundary earthquakes. *Earth Planets Space* 60(9):987–991
- Oglesby DD, Archuleta RJ (1997) A faulting model for the 1992 Petrolia earthquake: can extreme ground acceleration be a source effect? *J Geophys Res Solid Earth* 102(B6):11877–11897
- Olsen KB, Day SM, Bradley CR (2003) Estimation of Q for long-period (> 2 sec) waves in the Los Angeles basin. *Bull Seismol Soc Am* 93(2):627–638
- Olson AH, Orcutt JA, Frazier GA (1984) The discrete wavenumber/finite element method for synthetic seismograms. *Geophys J Int* 77(2):421–460
- Pitarka A (1999) 3D elastic finite-difference modeling of seismic motion using staggered grids with nonuniform spacing. *Bull Seismol Soc Am* 89(1):54–68
- Pu WC, Liang RJ, Dai FY, Huang B (2017) An analytical model for approximating pulse-like near-fault ground motions. *J Vib Shock* 36(4):208–213 (in Chinese)
- Ricker N (1943) Further developments in the wavelet theory of seismogram structure. *Bull Seismol Soc Am* 33(3):197–228
- Somerville P, Irikura K, Graves R, Sawada S, Wald D, Abrahamson N, Iwasaki Y, Kagawa T, Smith N, Kowada N (1999) Characterizing crustal earthquake slip models for the prediction of strong ground motion. *Seismol Res Lett* 70(1):59–80
- Wald DJ, Heaton TH, Hudnut KW (1996) The slip history of the 1994 Northridge, California, earthquake determined from strong-motion, teleseismic, GPS, and leveling data. *Bull Seismol Soc Am* 86(1B):S49–S70
- Wald DJ, Quitoriano V, Heaton TH, Kanamori H (1999) Relationships between peak ground acceleration, peak ground velocity, and modified Mercalli intensity in California. *Earthq Spectra* 15(3):557–564
- Wang Y (2015) The Ricker wavelet and the Lambert W function. *Geophys J Int* 200(1):111–115

- Xu LJ, Xie LL (2005) Characteristics of frequency content of near-fault ground motions during the Chi–Chi earthquake. *Acta Seismol Sin* 18(6):707–716 (**in Chinese**)
- Zeng Y, Anderson JG (1996) A composite source model of the 1994 Northridge earthquake using genetic algorithms. *Bull Seismol Soc Am* 86(1B):S71–S83
- Zhang WB, Yu XW (2010) Strong ground motion simulation of the 1999 Chi–Chi, Taiwan, earthquake. *J Earthq Eng Eng Vib* 30(3):1–11 (**in Chinese**)



# Thin interbed AVA inversion based on a fast algorithm for reflectivity

Zhen Yang<sup>1,2</sup> · Jun Lu<sup>3</sup>

Received: 9 January 2020 / Accepted: 18 May 2020 / Published online: 25 May 2020  
© Institute of Geophysics, Polish Academy of Sciences & Polish Academy of Sciences 2020

## Abstract

Zoeppritz equations form the theoretical basis of most existing amplitude variation with incident angle (AVA) inversion methods. Assuming that only primary reflections exist, that is, the multiples are fully suppressed and the transmission loss and geometric spreading are completely compensated for, Zoeppritz equations can be used to solve for the elastic parameters of strata effectively. However, for thin interbeds, conventional seismic data processing technologies cannot suppress the internal multiples effectively, nor can they compensate for the transmission loss accurately. Therefore, AVA inversion methods based on Zoeppritz equations or their approximations are not applicable to thin interbeds. In this study, we propose a prestack AVA inversion method based on a fast algorithm for reflectivity. The fast reflectivity method can compute the full-wave responses, including the reflection, transmission, mode conversion, and internal multiples, which is beneficial to the seismic inversion of thin interbeds. A further advantage of the fast reflectivity method is that the partial derivatives of the reflection coefficient with respect to the elastic parameters can be expressed as analytical solutions. Based on the Gauss–Newton method, we construct the objective function and model-updating formula considering sparse constraint, where the Jacobian matrix takes the form of an analytical solution, which can significantly accelerate the inversion convergence. We validate our inversion method using numerical examples and field seismic data. The inversion results demonstrate that the fast reflectivity-based inversion method is more effective for thin interbed models in which the wave-propagation effects, such as interval multiples, are difficult to eliminate.

**Keywords** Amplitude variation with incident angle · Inversion · Thin interbed · Fast algorithm · Reflectivity method

## Introduction

According to Aki and Richards (1980), the reflection and transmission coefficients depend on the incident angle and elastic parameters (the P- and S-wave velocities, as well as density). The technique of amplitude variation with offset (AVO) or amplitude variation with incident angle (AVA) inversion uses this dependency for the elastic parameter inversion. Owing to the complexity and nonlinearity of Zoeppritz equations, a considerable number of prestack AVA

inversion methods have been based on approximate solutions to Zoeppritz equations (Stewart 1990; Fatti et al. 1994; Larsen 1999; Jin 1999; Mahmoudian and Margrave 2004). Furthermore, to improve the inversion accuracy, many prestack AVA inversion methods based on the exact Zoeppritz equations method (EZM) have been developed (Tiğrek et al. 2005; Wang et al. 2011; Lu et al. 2015). Although the abovementioned inversion methods can achieve satisfactory results, they still exhibit several restrictions and limitations. Because the inversion methods based on Zoeppritz equations and their approximations use the assumption that only primary reflections are target wavefields, the wave-propagation effects, such as multiples, transmission loss, and geometric spreading, are completely eliminated. However, it is difficult for existing technology to deal with the above wave propagation effects without affecting the energy of primary reflections, particularly for thin interbedded formations.

Considering the wave propagation effects, the reflectivity method (RM) provides an elegant algorithm for computing full-wave reflection coefficients. Fuchs (1968) first proposed

---

✉ Jun Lu  
lujun615@163.com

<sup>1</sup> Key Laboratory of Earth and Planetary Physics, Institute of Geology and Geophysics, Chinese Academy of Sciences, Beijing 100029, China

<sup>2</sup> University of Chinese Academy of Sciences, Beijing 100049, China

<sup>3</sup> School of Geophysics and Information Technology, China University of Geosciences, Beijing 100083, China

RM for generating synthetic seismic seismograms. Later, the RM scheme was extended to include elastic transmission loss and time shift due to multiple layers on the top of reflecting medium (Fuchs and Müller 1971). They demonstrated that this improvement is necessary for practical RM applications, as only the reflections from the deeper parts of layered media are generally of interest, and the reflections from the deeper-layered media suffer from transmission losses. Kennett (1974) first proposed the recursive algorithm for the calculation of the total reflection and transmission coefficients for a stack of layers, in which the unconditional stability for all frequencies and slownesses was improved. The recursive algorithm basically solved the overflow problems in the calculation of exponential functions for high frequencies and slownesses (Kennett 1983, 2009). To improve the unsatisfactory computational efficiency of the RM, Phinney et al. (1987) proposed the fast reflectivity method (FRM) based on reorganization of the innermost loops of the Kennett RM (Kennett 1983). The FRM converts the Kennett RM into a vectorizable algorithm, which achieves a speed enhancement of approximately 20 times when implemented on an array processor (Phinney et al. 1987). Another advantage of the FRM is that the partial derivative of the reflection coefficient to the elastic parameters can be expressed as an analytical solution, which aids in improving the inversion speed and accuracy.

Under the assumption of the locally one-dimensional (1D) model, the RM has been used extensively in layered stratigraphic inversion. Sen and Roy (2003) compared the characteristics of the Kennett RM and FRM, and then designed a regularized Gauss–Newton-type algorithm for prestack waveform inversion by rearranging the recursion formula in the Kennett RM. However, this inversion was performed in the intercept time and slowness domain. To prevent the aliasing problem, the slowness must be adequately sampled, but the calculation time will increase exponentially with the slowness samples (Mallick and Frazer 1987). Liu et al. (2016) developed a Bayesian inversion methodology for the P-wave AVO inversion method based on the FRM. They implemented a modification to convert the inversion from the intercept time and slowness domain into the angle gather domain. Liu et al. (2018) proposed the FRM-based nonlinear multicomponent prestack AVA joint inversion method using the non-dominated sorting genetic algorithm.

In this study, we propose an FRM-based prestack AVA inversion method for thin interbed strata. Based on the theories of the Kennett RM and FRM, we establish the objective function by means of the least-squares approach. To estimate the elastic parameters (P- and S-wave velocity, as well as density), the objective function is achieved by the minimized difference between the simulated and observed data in the angle domain. Moreover, in our inversion theory, the Jacobian matrix is expressed in the form of an analytical

solution, and we describe the derivation process in detail in the Appendix. Eventually, we implement the FRM-based inversion on the thin interbed model and field data, and then compare with the inversion method based on the EZM. The inversion results demonstrate that the FRM-based inversion approach is better than the EZM-based inversion technique in terms of accuracy and continuity to thin interbeds.

## Theory

### Kennett reflectivity method

Considering two consecutive regions AB and BC (Fig. 1), Kennett (1983) derived the following equations using the recursive method (Kennett 2009; p. 104), which can calculate the overall reflection and transmission matrices of region AC when those of regions AB and BC are known:

$$\begin{cases} \mathbf{R}_D^{AC} = \mathbf{R}_D^{AB} + \mathbf{T}_U^{AB} \mathbf{R}_D^{BC} [\mathbf{I} - \mathbf{R}_U^{AB} \mathbf{R}_D^{BC}]^{-1} \mathbf{T}_D^{AB}, \\ \mathbf{T}_D^{AC} = \mathbf{T}_D^{BC} [\mathbf{I} - \mathbf{R}_U^{AB} \mathbf{R}_D^{BC}]^{-1} \mathbf{T}_D^{AB}, \\ \mathbf{R}_U^{AC} = \mathbf{R}_U^{AB} + \mathbf{T}_D^{BC} \mathbf{R}_U^{AB} [\mathbf{I} - \mathbf{R}_D^{BC} \mathbf{R}_U^{AB}]^{-1} \mathbf{T}_U^{BC}, \\ \mathbf{T}_U^{AC} = \mathbf{T}_U^{AB} [\mathbf{I} - \mathbf{R}_D^{BC} \mathbf{R}_U^{AB}]^{-1} \mathbf{T}_U^{BC}, \end{cases} \quad (1)$$

where the superscripts denote the regions; the subscripts D and U denote the downgoing and upgoing waves, respectively; and  $\mathbf{R}$  and  $\mathbf{T}$  are the reflection and transmission matrices.

Equation (1) presents the most fundamental development of the reflectivity formulations. This unconditionally stable algorithm can compute the full-wave response, including the interlayer wave propagation and interaction terms. Furthermore, Fig. 1 clearly indicates that Eq. (1) includes all of the internal multiples and mode-converted waves. It is also possible to compute certain selected modes under the

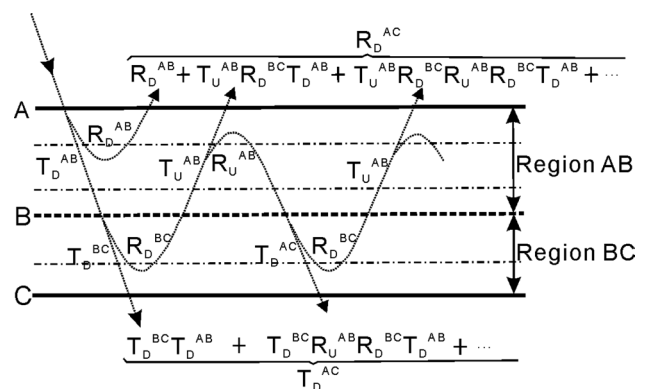


Fig. 1 Schematic representation of first several terms of expansion of addition rules for reflection and transmission matrices, indicating interactions with regions ‘AB’ and ‘BC’ (Kennett 2009)

reflectivity formulation. Detailed descriptions of the theory and discussions can be found in the book by Kennett (2009).

### Fast reflectivity method

In applications, the conventional RM described above is unsatisfactorily inefficient. Phinney et al. (1987) proposed a speed-up algorithm based on the reorganization of the inside loops of the conventional RM to permit vectorization. The vectorized procedure by itself is faster than the previously described procedures, simply because practically all of the redundancy has been eliminated. According to Phinney et al. (1987), the overall reflection coefficient of the PP-wave in the frequency domain is

$$R(p, \omega) = \frac{v_{04}}{v_{01}} \tag{2}$$

where  $p$  is the horizontal slowness;  $\omega$  is the angular frequency; and  $v_{01}$  and  $v_{04}$  are the first and fourth elements of the vector  $v_0$ , respectively. Moreover,  $v_0$  is a vector with six elements, as defined by Phinney et al. (1987) in the frequency–slowness domain ( $\omega$ – $p$ ) based on the Haskell–Dunkin minor matrix:

$$v_0 = [\Delta \ -R_{PS}\Delta \ -R_{SS}\Delta \ R_{PP}\Delta \ R_{SP}\Delta \ \det R\Delta]^T, \tag{3}$$

where  $\Delta$  is a scaling factor;  $R_{PP}$ ,  $R_{PS}$ ,  $R_{SP}$ , and  $R_{SS}$  are the reflection coefficients of the PP-, PS-, SP-, and SS-waves, respectively; and  $\det R$  is the determinant of the coefficients.

Supposing a stack of  $N$  isotropic horizontal layers, we can obtain  $v_0$  starting from  $v_N$  via a sequence of matrix multiplications:

$$v_0 = Q_0 Q_1 \dots Q_n \dots Q_{N-1} v_N, \quad n \in [0, N - 1], \tag{4}$$

where  $v_N$  is the initial six-element vector

$$v_N = [1 \ 0 \ 0 \ 0 \ 0 \ 0]^T, \tag{5}$$

and  $Q_n$  is the wave propagator matrix of the  $n$ th layer. Each  $Q_n$  is the product of an interface matrix  $F_n$  and a layer-crossing matrix  $E_n$ , as follows:

$$Q_n = E_n F_n. \tag{6}$$

The specific forms of  $E_n$  and  $F_n$  are as follows:

$$E_n = \text{diag} [ e^{-i\omega h_n(q_n^p + q_n^s)} \ 1 \ e^{-i\omega h_n(q_n^p - q_n^s)} \ e^{-i\omega h_n(q_n^p - q_n^s)} \ 1 \ e^{-i\omega h_n(q_n^p + q_n^s)} ], \tag{7}$$

$$F_n = T_n^{-1} T_{n+1}, \tag{8}$$

where  $h_n$  is the layer thickness of the  $n$ th layer;  $q_n^p$  and  $q_n^s$  are the vertical slownesses of P and S;  $T_n$  is the  $6 \times 6$  delta matrix for the  $n$ th layer; and  $T_n^{-1}$  is the inverse matrix of  $T_n$ . Details regarding  $T_n$  and  $T_n^{-1}$  are presented in “Appendix 1”.

For thick layers, multiples can be suppressed by the predictive deconvolution method. However, for thin-bed and thin-interbed layers, the primary waves always interfere with internal multiples, which are difficult to eliminate using the existing technologies without destroying the amplitudes of primary waves. Therefore, internal multiples must be considered in the inversion for a single thin bed or thin interbed due to their significant impact on the overall reflection coefficients. Besides, because there is no layer thickness limitation, both the Kennett reflectivity method and fast reflectivity method can be applied to thin-bed and interbed strata.

### P-wave AVA inversion

#### Gauss–Newton algorithm

We use the Gauss–Newton algorithm to solve nonlinear least-squares problems in our AVA inversion. The inversion, which uses the FRM to simulate synthetic data, is achieved by the minimized differences between the simulated and observed data, so as to estimate the elastic parameters (P-wave velocity  $\alpha$ , S-wave velocity  $\beta$ , and density  $\rho$ ). The objective function is formulated as follows:

$$Q(\mathbf{M}) = \|\Phi - \Phi^{\text{obs}}\|^2, \tag{9}$$

where  $\mathbf{M} = (\alpha, \beta, \rho)$  denotes the model parameter vector in the target time window, while  $\Phi^{\text{obs}}$  and  $\Phi$  are the observed and synthetic data, respectively.

It should be noted that  $\Phi^{\text{obs}}$  and  $\Phi$  are angular gathers in the time domain, while the reflection coefficient calculated by Eq. (2) is in the frequency–slowness domain; therefore, domain conversion is required. The entire domain transformation process can be divided into the following two steps.

#### a. Slowness to incident angle

The incident angle of each trace in the angle gather is fixed. Thus, under the assumption of horizontal layers, we can obtain the corresponding slowness of the incident angle for each layer according to Snell’s law:

$$p_n = \frac{\sin \theta}{\alpha_n}, \tag{10}$$

where  $\theta$  is the incident angle and the subscript  $n$  indicates the  $n$ th layer. Then, the reflection coefficient in the slowness domain can be converted into the reflection coefficient in the time domain, which is  $R(p, \omega)$  to  $R(\theta, \omega)$ .

#### b. Frequency domain to time domain

We implement the conversion from the frequency domain to the time domain by means of an inverse Fourier transform. The outcome is the reflection coefficient

datasets in the intercept time and angle ( $\tau$ - $\theta$ ) domain (Fryer 1980):

$$R(\tau, \theta) = \frac{1}{2\pi} \int_{-\infty}^{\infty} R(\omega, \theta) e^{i\omega\tau} d\omega. \tag{11}$$

Following the domain conversion, the synthetic seismogram can be obtained by the convolution of the reflection coefficients and wavelets:

$$\Phi = \mathbf{W} * \mathbf{R}, \tag{12}$$

where  $\mathbf{W}$  denotes the wavelets of the PP-waves at different incident angles and  $\mathbf{R}$  is the matrix of the PP-wave reflection coefficients.

**Model update**

Under approximately ideal conditions, given an initial guess model  $\mathbf{M}_0$  close to the true model, we can obtain the model update matrix using the Gauss–Newton method (Tarantola 1986; Sheen et al. 2006; Lu et al. 2015, 2017):

$$\Delta\mathbf{M} = [\mathbf{J}^T \mathbf{J}]^{-1} \mathbf{J}^T (\Phi - \Phi^{obs}), \tag{13}$$

where  $\Delta\mathbf{M} = (\Delta\alpha, \Delta\beta, \Delta\rho)$  is the update of the initial model and  $\mathbf{J}$  is the Jacobian matrix:

$$\mathbf{J} = \frac{\partial(\mathbf{W} * \mathbf{R})}{\partial\mathbf{M}} = \mathbf{W} * \frac{\partial\mathbf{R}}{\partial\mathbf{M}}. \tag{14}$$

Following the damped least-squares method (Levenberg 1944; Marquardt 1963), the regularized Gauss–Newton formula is:

$$\Delta\mathbf{M} = [\mathbf{H} + \lambda\mathbf{I}]^{-1} \mathbf{J}^T (\Phi - \Phi^{obs}), \tag{15}$$

where  $\lambda$  is a scalar,  $\mathbf{I}$  is the identity matrix, and  $\mathbf{H} = \mathbf{J}^T \mathbf{J}$  is the Hessian matrix. The details of  $\lambda$  are described by Paige and Saunders (1982).

According to Eq. 15, we build an objective function considering sparse constraint (Yuan et al. 2019; Luo et al. 2018) for our AVO inversion as

$$Q = \left\| (\mathbf{H} + \lambda\mathbf{I})\Delta\mathbf{M} - \mathbf{J}^T (\Phi - \Phi^{obs}) \right\|^2 + k^2 \|\mathbf{M}_{inv} - \mathbf{M}_{true}\|, \tag{16}$$

where  $k$  is the weight for the sparse constraint, and the method of choosing  $k$  can be found in the paper (Chen et al. 2001). Letting  $Q$  reach to a minimum, we can approximately derive

$$(\mathbf{H} + \lambda\mathbf{I})\Delta\mathbf{M} - \mathbf{J}^T (\Phi - \Phi^{obs}) \approx 0, \tag{17}$$

$$k^2 (\mathbf{M}_{inv} - \mathbf{M}_{true}) \approx 0. \tag{18}$$

Under the approximate ideal conditions, given an initial guess model  $\mathbf{M}_0$  close to the true model, we derive the approximate model update matrix as

$$\Delta\mathbf{M} \approx (\mathbf{H} + \lambda\mathbf{I} + k^2\mathbf{I})^{-1} [\mathbf{J}^T (\Phi - \Phi^{obs}) + k^2(\mathbf{M}_0 - \mathbf{M}_{true})]. \tag{19}$$

In the subsequent iteration,  $\Delta\mathbf{M}$  is used to update the initial  $\mathbf{M}_0$ , and the iteration stops when  $\Phi^{obs}$  is close to  $\Phi$  for a given accuracy.

**Jacobian matrix**

According to the definition of the reflection coefficients in Eq. (11), the analytical Jacobian matrix of the PP-wave is

$$\mathbf{J} = \mathbf{W} * \frac{1}{2\pi} \int_{-\infty}^{\infty} \frac{\partial R(\omega, \theta)}{\partial \mathbf{M}_n} e^{i\omega\tau} d\omega, \tag{20}$$

where

$$\frac{\partial R(\omega, \theta)}{\partial \mathbf{M}_n} = \frac{v_{01} \frac{\partial v_{04}}{\partial \mathbf{M}_n} - v_{04} \frac{\partial v_{01}}{\partial \mathbf{M}_n}}{(v_{01})^2}. \tag{21}$$

From Eqs. (6) and (8), it can be found that only  $\mathbf{Q}_{n-1}$  and  $\mathbf{Q}_n$  are related to the parameters  $\mathbf{M}_n$ . Therefore, the partial derivative of  $v_0$  is

$$\frac{\partial v_0}{\partial \mathbf{M}_n} = \mathbf{Q}_0 \mathbf{Q}_1 \cdots \frac{\partial(\mathbf{Q}_{n-1} \mathbf{Q}_n)}{\partial \mathbf{M}_n} \cdots \mathbf{Q}_{N-1} v_N, \tag{22}$$

where

$$\frac{\partial(\mathbf{Q}_{n-1} \mathbf{Q}_n)}{\partial \mathbf{M}_n} = \frac{\partial \mathbf{Q}_{n-1}}{\partial \mathbf{M}_n} \mathbf{Q}_n + \mathbf{Q}_{n-1} \frac{\partial \mathbf{Q}_n}{\partial \mathbf{M}_n}, \tag{23}$$

$$\frac{\partial \mathbf{Q}_n}{\partial \mathbf{M}_n} = \frac{\partial E_n}{\partial \mathbf{M}_n} \mathbf{F}_n + \mathbf{E}_n \frac{\partial \mathbf{F}_n}{\partial \mathbf{M}_n}, \tag{24}$$

in which  $\mathbf{M}_n$  denotes  $\alpha, \beta,$  and  $\rho$  of the  $n$ th layer.

The partial derivative of  $\mathbf{F}_n$  is

$$\frac{\partial \mathbf{F}_n}{\partial \mathbf{M}_n} = \frac{\partial \mathbf{T}_n^{-1}}{\partial \mathbf{M}_n} \mathbf{T}_{n+1} + \mathbf{T}_n^{-1} \frac{\partial \mathbf{T}_{n+1}}{\partial \mathbf{M}_n}, \tag{25}$$

where

$$\frac{\partial \mathbf{T}_{n+1}}{\partial \mathbf{M}_n} = \begin{cases} \frac{\partial \mathbf{T}_{n+1}}{\partial \alpha}, \\ \frac{\partial \mathbf{T}_{n+1}}{\partial \beta}, \\ \frac{\partial \mathbf{T}_{n+1}}{\partial \rho}. \end{cases} \tag{26}$$

The partial derivative of  $\mathbf{T}_n$  for the parameters  $\mathbf{M}_n$  is presented in ‘‘Appendix 2’’. We can also analytically calculate the partial derivative of  $\mathbf{E}_n$  for  $\mathbf{M}_n$ , which is shown in ‘‘Appendix 2’’.

Similar to the accelerated method presented by Phinney et al. (1987), we extend  $\frac{\partial \mathbf{E}_n}{\partial \mathbf{M}_n}$  to  $\hat{\mathbf{E}}_{\mathbf{M}_n}$  by including all frequencies, which is the extension process from a specific frequency to all frequencies for the  $n$ th layer. To optimize the computation of  $\hat{\mathbf{E}}_{\mathbf{M}_n}$ , we generate the values for  $\hat{\mathbf{E}}_{\mathbf{M}_n}$  recursively, assigning the first six elements (zero frequency) as the starting values. Consequently,  $\hat{\mathbf{E}}_{\mathbf{M}_n}$  is computed using the increment vector  $\mathbf{g}'_n$  to generate the higher-frequency values of  $\hat{\mathbf{E}}_{\mathbf{M}_n}$  by means of complex multiplication.

$$\begin{bmatrix} \hat{E}_{\mathbf{M}_n} \\ \hat{E}_{\mathbf{M}_{n+1}} \\ \vdots \\ \hat{E}_{\mathbf{M}_{n+5}} \end{bmatrix} = \begin{bmatrix} g'_n(1) \\ g'_n(2) \\ \vdots \\ g'_n(6) \end{bmatrix} \begin{bmatrix} \hat{E}_{\mathbf{M}_{n-6}} \\ \hat{E}_{\mathbf{M}_{n-5}} \\ \vdots \\ \hat{E}_{\mathbf{M}_{n-1}} \end{bmatrix}, m = 7, 13, \dots, 6 \times (\text{NF} - 1) + 1, \tag{27}$$

where NF is the number of frequency.

The increment vector  $\mathbf{g}'_n$  is defined as

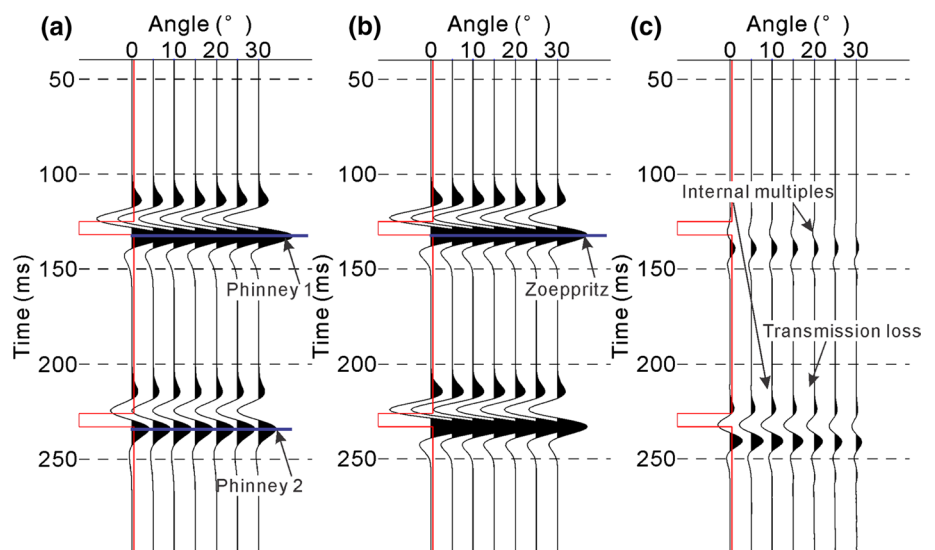
$$\mathbf{g}'_n = \begin{bmatrix} -\frac{\omega_m}{\omega_{m-6} s_n^+} & 0 & \frac{\omega_m}{\omega_{m-6} s_n^+} & \frac{\omega_m s_n^-}{\omega_{m-6}} & 0 & \frac{\omega_m s_n^+}{\omega_{m-6}} \\ \omega_{m-6} s_n^+ & \omega_{m-6} s_n^+ & \omega_{m-6} s_n^+ & \omega_{m-6} s_n^+ & \omega_{m-6} s_n^+ & \omega_{m-6} s_n^+ \end{bmatrix}, m = 7, 13, \dots, 6 \times (\text{NF} - 1) + 1. \tag{28}$$

where

**Table 1** Parameters of the theoretical model

Layer	$V_p$ (m/s)	$V_s$ (m/s)	Density (g/cm <sup>3</sup> )	Thickness (m)
1	3200	1816	2.5	200
2	2200	1300	1.5	8
3	3200	1816	2.5	150
4	2200	1300	1.5	8
5	3200	1816	2.5	200

**Fig. 2** Comparison of PP-wave synthetic angle gathers based on FRM and EZM: **a** angle gather based on FRM, **b** angle gather based on EZM, and **c** difference between angle gathers based on FRM and EZM. The red curve indicates the P-wave velocity



$$\begin{cases} s_n^+ = e^{i\Delta\omega\Delta_n(q^p+q^s)}. \\ s_n^- = e^{i\Delta\omega\Delta_n(q^p-q^s)}. \end{cases} \tag{29}$$

When we extend  $\frac{\partial \mathbf{E}_n}{\partial \mathbf{M}_n}$  to  $\hat{\mathbf{E}}_{\mathbf{M}_n}$ , we need suitable extensions of  $\mathbf{F}_n$  and  $\mathbf{v}_0$  using the same method as that in Phinney et al. (1987).

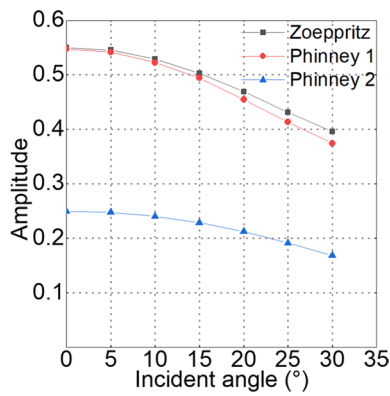
### Synthetic data test

#### Comparison of FRM and EZM

To verify the effectiveness of the FRM in AVA modeling and inversion, we used a horizontally layered model (as indicated in Table 1). The PP-wave synthetic angle gather was generated by the FRM and EZM with a 40 Hz Ricker wavelet. The detailed simulation steps are as follows.

1. Calculate zero-offset travel time.
2. Given the incident angle, calculate the corresponding reflection coefficient.
3. Calculated the angle gather by the convolution of the reflection coefficient and wavelet matrices.

Figure 2 shows that compared with the gathers based on the EZM, those based on the FRM contain more complete information (such as the multiples illustrated in Fig. 2c). This is because the synthetic seismogram using the EZM only reflects the primary reflection amplitudes, without



**Fig. 3** Comparison of AVA curves based on FRM and EZM. Phinney 1 and Phinney 2 are two interfaces shown in Fig. 2a, and Zoeppritz is the interface shown in Fig. 2b

**Table 2** Parameters of thin interbed model

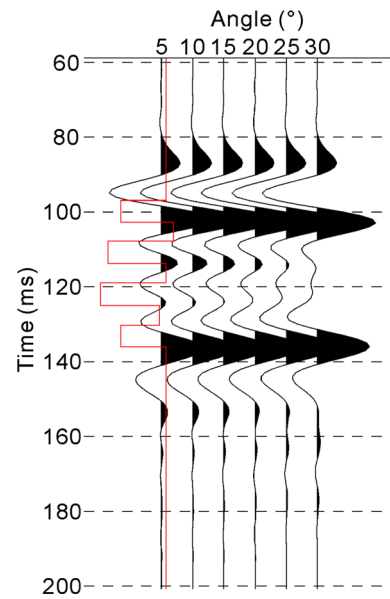
$V_p$ (km/s)	$V_s$ (km/s)	Density ( $g/cm^3$ )	Thickness (m)
3.094	1.515	2.4	150
2.781	1.665	2.08	8
3.146	1.554	2.41	8
2.694	1.206	2.3	8
3.094	1.515	2.4	8
2.643	1.167	2.29	8
3.048	1.595	2.23	8
2.781	1.665	2.08	8
3.094	1.515	2.4	150

containing wave-propagation effects. These multiples may cause errors in interpretation and inversion.

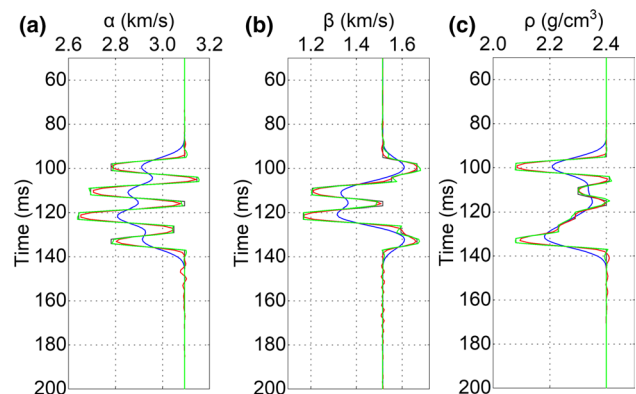
To display the difference between the FRM and EZM further, we extracted the AVA curves of two interfaces (illustrated in Fig. 3) with the same elastic parameters on either side of the interface (along the blue line in Fig. 2). Owing to the influence of transmission loss, the amplitudes of the three interfaces using the reflectivity method were gradually reduced with the depth and were smaller than those obtained using the EZM. If the transmission loss and internal multiples are not properly corrected in seismic data processing procedures, the synthetics obtained by the FRM can be more effectively matched with the field data.

**Thin interbed model test**

To test the inversion method, we selected a 1D thin interbed model, in which the thicknesses of all single thin layers were set to 8 m (as indicated in Table 2). The corresponding PP synthetic AVA gather was generated through the convolution of the reflectivity derived from the FRM and the Ricker wavelets, with dominant frequencies of 40 Hz.



**Fig. 4** Synthetic PP seismograms based on FRM. The red curve denotes the P-wave velocity

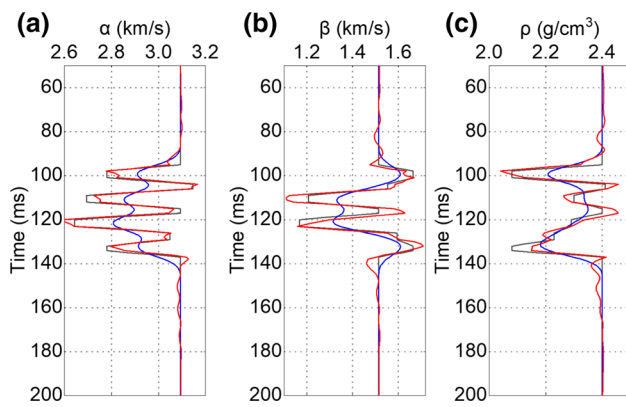


**Fig. 5** Inversion results of **a** P-wave velocity, **b** S-wave velocity, and **c** density by FRM-based method. The blue, black, red, and green curves denote the initial model, true model, inverted result, and inverted result for sparse constraint, respectively

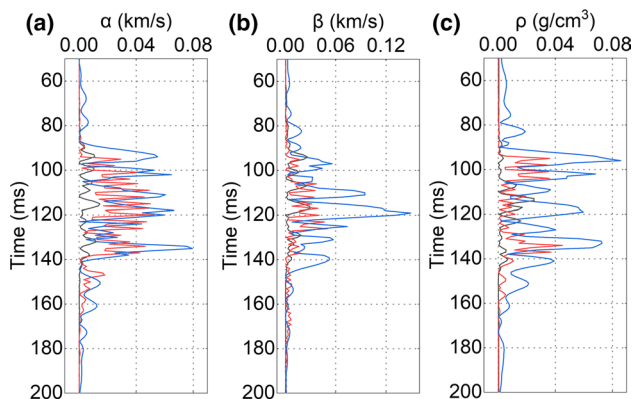
**Inversion of synthetic data without noise**

The AVA inversions based on the FRM and EZM were tested. Both of the inversion methods used the same input data (as illustrated in Fig. 4) and the same initial model (the blue line in Fig. 5). The FRM- and EZM-based inversion results of the P- and S-wave velocities, as well as density, are presented in Figs. 5 and 6, respectively. It can be observed that the inversion parameters based on the FRM basically matched the true values. However, the inversion based on the EZM regarded internal multiples as primary reflections, leading to false images and instability. Considering the sparse constraint during the FRM-based inversion, we derived the improved





**Fig. 6** Inversion results of **a** P-wave velocity, **b** S-wave velocity, and **c** density by EZM-based method. The blue, black, and red curves denote the initial model, true model, and inverted result, respectively

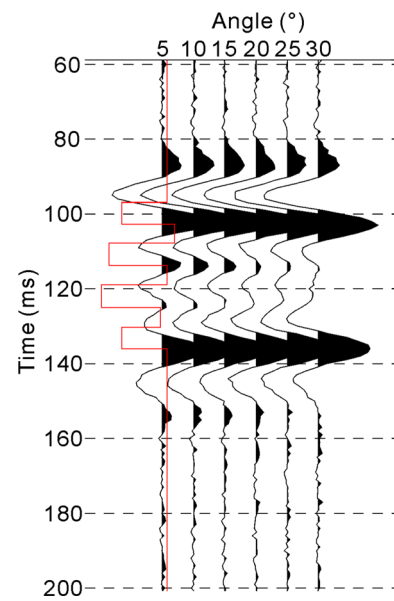


**Fig. 7** Absolute values of the differences between the true curves and inversion results of different methods. The black, red, and blue curves denote the differences between the true curves and FRM-based inversion results considering sparse constraint, between the true curves and FRM-based inversion results, and between the true curves and EZM-based inversion results, respectively

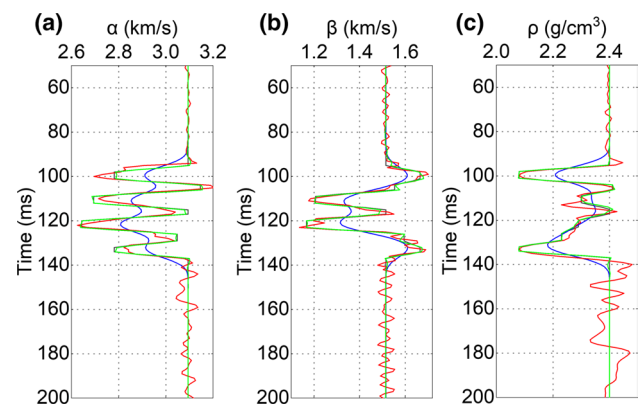
inversion results shown in the green curves in Fig. 5. In order to quantitatively describe the inversion results of different inversion methods, we calculated the differences between the inversion results and true values, as shown in Fig. 7. It can be found that the inversion results considering the sparse constraint are closer to the true curves.

### Inversion of synthetic data with noise added

As illustrated in Fig. 8, the robustness of the inversion method was then tested on the synthetic PP angle gather with a 15% level of random noise added. The FRM- and ZEM-based inversion results using the same initial model are presented in Figs. 9 and 10, respectively. It can be observed that under the influence of random noise in the angle gather, the inversion results from both inversion methods were noisy. However,

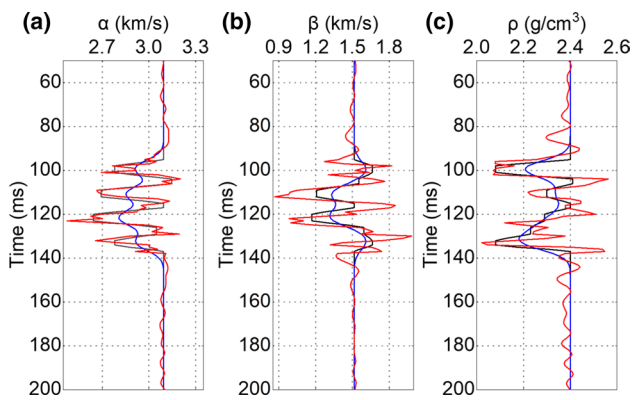


**Fig. 8** Synthetic PP seismograms based on FRM with 15% noise level. The gray curve denotes the P-wave velocity



**Fig. 9** Inversion results of **a** P-wave velocity, **b** S-wave velocity, and **c** density by FRM-based approach under noisy conditions. The blue, black, red, and green curves denote the initial model, true model, inverted result, and inverted result for sparse constraint, respectively

the inversion results at the thin interbedded layers could still reflect the information of the true model parameters. Moreover, the FRM-based inversion results were closer to the true values, indicating that the FRM exhibited stronger robustness. As shown in the green curves in Fig. 9, if we consider the sparse constraint during the inversion, we can hardly find the effect of random noise on the inversion results. Compared with the green curves in Fig. 5, the inversion results from the noisy condition are almost unchanged. In order to quantitatively describe the inversion results of different inversion methods, we calculated the differences between the inversion results and true values, as shown in Fig. 11. We can still find that even with the noise, the inversion results considering the

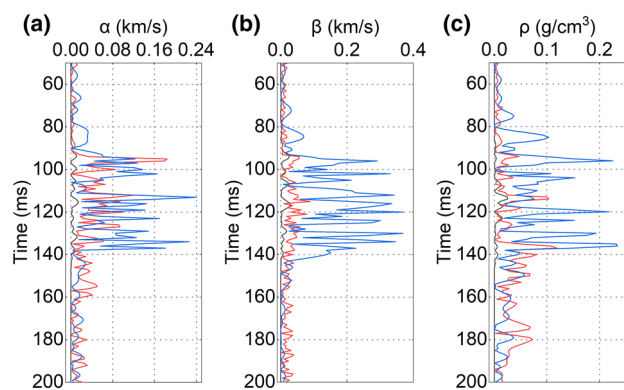


**Fig. 10** Inversion results of **a** P-wave velocity, **b** S-wave velocity, and **c** density by EYM-based approach under noisy conditions. The blue, black, and red curves denote the initial model, true model, and inverted result, respectively

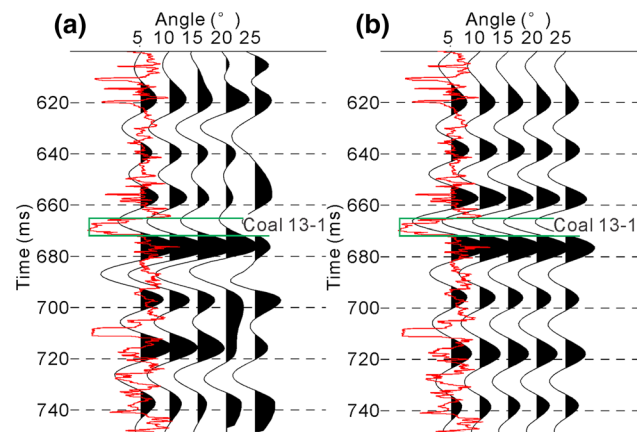
sparse constraint fit better with the true values. Therefore, our method has a certain ability to resist noise.

### Application to field data

Our inversion method was also applied to the processed PP field AVA datasets (Fig. 12a) from Guqiao mine, located in the Huainan coalfield on the southern margin of the North China plate. The 3D seismic data were acquired in 2006, which covered an area of 2.56 km<sup>2</sup> with a bin size of 10 × 10 m. The coal-bearing strata dip gently with an angle less than 5°, and the structure is relatively simple with few faults. The 3D acquisition offsets ranged from 0 to 1210 m, which led to 0°–35° incident angles on coal seam 13–1. In



**Fig. 11** Absolute values of the differences between the true curves and inversion results of different methods under noisy conditions. The black, red, and blue curves denote the differences between the true curves and FRM-based inversion results considering sparse constraint, between the true curves and FRM-based inversion results, and between the true curves and EYM-based inversion results, respectively



**Fig. 12** Angle gathers for field data and synthetic data: **a** field angle gather data and **b** synthetic angle data. The red curve denotes the density logs, while the green frame represents the location of coal seam 13–1

this application, we chose one 2D line across the well to show the inversion effect of the proposed method. The major stratigraphic units illustrated in Fig. 12a are the coal measure strata, where commercial coal beds have been developed. The coal seams in the angle gather can be identified with the aid of log curves (the red curve in Fig. 12), where the strata with a low density and velocity of 665–675 ms is coal seam 13–1. Moreover, the lithology of the surrounding rocks is dominated by sand and shale, exhibiting a thin interbed structure. We only acquired the acoustic and density logs at the well location. Then, for the area consisting of sand and mudstone strata, we adopted an empirical correlation between the P- and S-wave velocities to convert the acoustic logs into the S-wave logs (Lu et al. 2016).

$$V_S = 0.433V_P + 430.9 \tag{30}$$

For the transformation of the S-wave log into the coal seam, we adopted the empirical correlation provided by Wang et al. (2016):

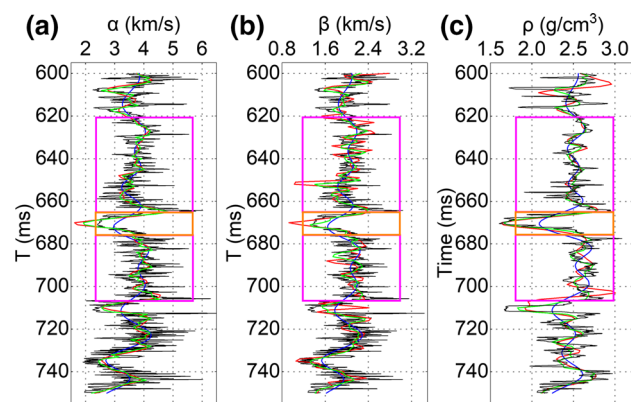
$$V_S = 0.5208V_P + 110.67. \tag{31}$$

Before implementing our inversion method, firstly, we had to ensure that the datasets were appropriately processed. Since the FRM-based inversion considers the transmission losses and internal multiples, the transmission loss compensation and internal multiples suppression are not adopted in the data processing. Secondly, with the aid of synthetic seismograms calculated from the well logs (Fig. 12b), we could calibrate the logs (depth domain) and PP events reflected from each geological interface (time domain). The correlation coefficient between the actual and synthetic PP AVA datasets was 0.74. Thirdly, using the well logs and

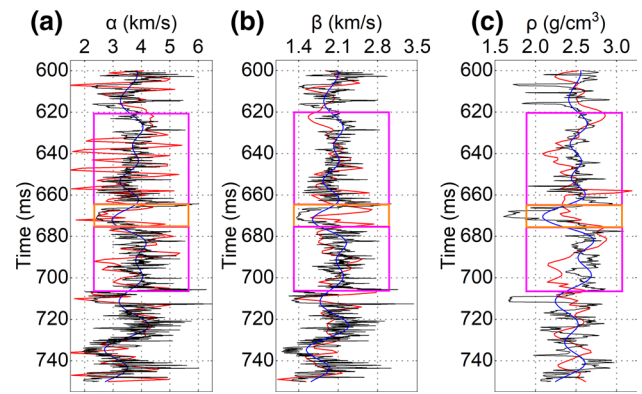
correlated angle gather, we estimated the wavelets, which were angle independent. The final step was to calculate the initial model, which was obtained by band-pass filtering, in which the range was 0–60 Hz.

Figures 13 and 14 display the inversion results using the two inversion methods at the well position. Figure 15 shows the difference between the inversion results and true values. The inversion results considering the sparse constraint have higher resolutions. Besides, compared with the EZM-based inversion result, the majority of variation trends of the FRM-based inversion results closely matched the well logs. For example, the inversion results based on the FRM could accurately describe the position of coal seam 13–1, and the inversion results of the P- and S-wave velocities, as well as the density, were strongly matched with the logging data. Moreover, we could distinguish more thin layers from the FRM-based inversion results than from the EZM-based inversion results. Here, we take the P-wave velocity inversion results as an example. In Figs. 13a and 14a, the strata in the pink frame are thin interbedded strata. The inversion results illustrated in Fig. 13a can reflect the variation trend of the strata, while those in Fig. 14a exhibit significant fluctuations.

Figure 16 shows the FRM-based inversion results of the 2D line across the well, where the black colors represent coal seams, which are calibrated by the well log. The wavelets used for inversion were extracted from the angle gathers at all incident angles. Although the inversion results are of lower frequency compared with the well logs due to the limitation of seismic resolution, the FRM-based inversion method can produce high-resolution results, which is very helpful for identifying thin layers. According to the well tops and seismic horizons, we picked the top and bottom interfaces of coal seam 13–1 (green lines in Fig. 16). It is



**Fig. 13** FRM-based inversion results for field data considering sparse constraint: **a** P-wave velocity, **b** S-wave velocity, and **c** density. The blue, black, red, and green curves represent the initial model, true model, inversion result, inversion result considering sparse constraint, respectively. The orange frame represents the location of coal seam 13–1, and the pink frames represent the thin interbedded strata

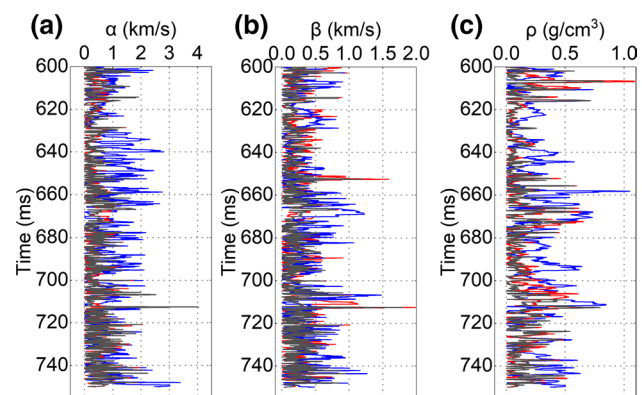


**Fig. 14** EZM-based inversion results for field data: **a** P-wave velocity, **b** S-wave velocity, and **c** density. The blue, black, and red curves denote the initial model, true model, and inverted result, respectively. The orange frame represents the location of coal seam 13–1, and the pink frames represent the thin interbedded strata

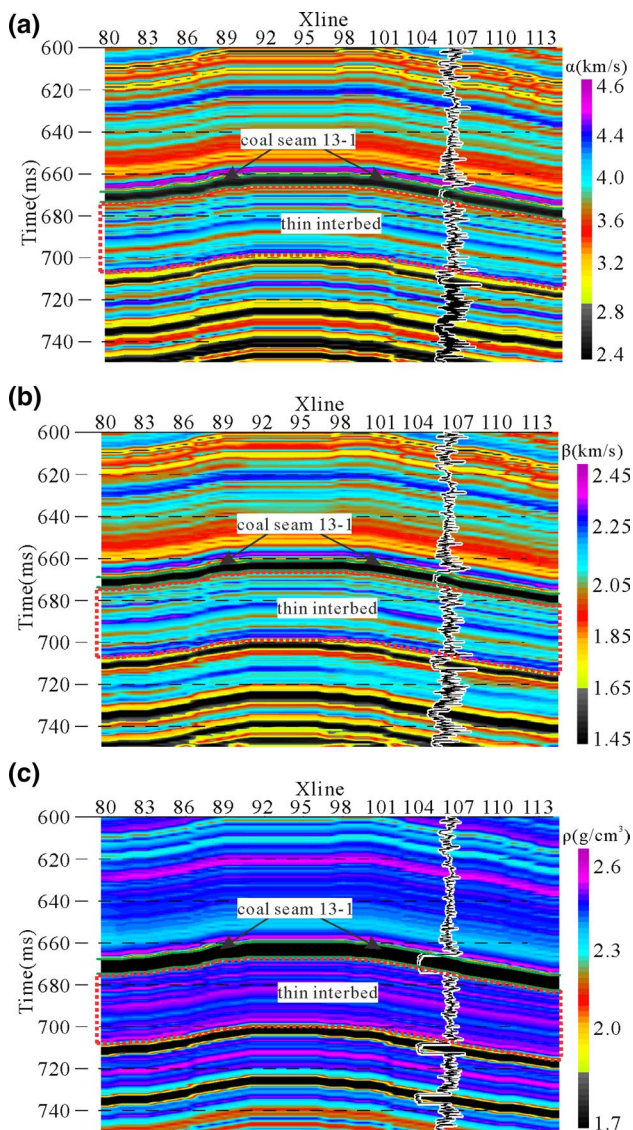
seen that the thickness of coal seam 13–1 is close to that indicated by the well log. Moreover, we can clearly find the thin interbedded strata under coal seam 13–1 in the inversion sections (the red frames in Fig. 16).

## Discussion

In seismic data processing, it is difficult to compensate for transmission loss completely using existing techniques. Although certain scholars (Xu et al. 1998; Zhang et al. 2003) attempted to compensate for transmission losses using post- or prestack migration, transmission loss compensation has not yet been sufficiently accurate (Deng and McMechan,



**Fig. 15** Absolute values of the differences between the true curves and inversion results of different methods for field data. The black, red, and blue curves denote the difference between the true curves and FRM-based inversion results considering sparse constraint, between the true curves and FRM-based inversion results, and between the true curves and EZM-based inversion results, respectively



**Fig. 16** FRM-based inversion sections for field data considering sparse constraint: **a** P-wave velocity, **b** S-wave velocity, and **c** density. The inserted black curves are the well logs corresponding to the inverted parameter types. The green lines indicate the top and bottom interfaces of coal seam 13–1. The red frames mark the thin interbeded strata

2007). However, the FRM takes the transmission loss into account during the reflectivity calculation, which can reduce the difficulty of prestack processing to a certain extent.

Numerous methods for weakening multiple amplitudes in seismic data are available (Weglein et al. 2011). However, the suppression of multiple waves, such as internal multiples, remains a challenge in data processing. Furthermore, the internal multiples are often mixed with the primary reflections of the surrounding rock. Therefore, the suppression of inter-layer multiples often destroys the amplitude of the primary reflections to an extent. In this case, the FRM considers the

multiples when calculating the reflection coefficients, which allows us to invert the target data without multiple attenuation.

When applying the FRM-based inversion method, the initial models are low-frequency models, which are obtained from the smoothed logs and must be able to reflect the strata trend approximately. As the frequency of the initial model increases, the rate of inversion convergence increases. Besides, if we adopt an initial model without the strata trend, such as the random model or linear model, the inversion will not achieve convergence.

The FRM-based inversion approach has stronger anti-noise robustness than the EZM-based inversion technique. However, if the noise is as strong as the reflection amplitude of a thin layer, the inversion cannot distinguish between the thin layer reflection and noise, which will cause the inversion results to fluctuate.

## Conclusions

According to the comparison of the synthetic angle gathers based on the FRM and EZM, the FRM takes into account the transmission loss and multiple waves when calculating the reflectivity. Moreover, we discussed the FRM-based inversion theory using the Gauss–Newton algorithm and tested the inversion method on thin interbeded model data as well as field seismic data. It should be emphasized that the processing procedure prior to inversion should include surface-related multiple attenuations and geometric spreading compensation, but it should exclude transmission loss compensation and internal multiples attenuation.

Based on the least squares approach, we considered the sparse constraint in the inversion. The model tests show that the proposed inversion method has stronger anti-noise robustness when considering the sparse constraint. Application to the field data inversion results demonstrates that our inversion method can effectively improve the resolution of thin interbeded.

**Acknowledgements** The authors are very grateful to the MWMC Group for processing the seismic data. We would also like to express thanks for the sponsorship of the National Natural Science Foundation of China (Nos 41574126 and U1910205).

**Author contributions** Each author has contributed to the present paper. Jun Lu conceived the idea of this research. Zhen Yang and Jun Lu designed and programmed the codes. Zhen Yang performed the simulation tests. Zhen Yang and Jun Lu applied the method to the field data and analyzed the inversion results. The paper was written by all the authors.

## Compliance with ethical standards

**Conflict of interest** The authors declare that they have no conflict of interest.

### Appendix 1

$T_n$  is the  $6 \times 6$  delta matrix for the  $n$ th layer. The elements of  $T_n$  are formed from the Dunkin matrix, which is frequency independent. The following is a list of the 16 independent elements of the matrix  $T_n$ :

$$\begin{aligned}
 t_{11} &= -(p^2 + q^p q^s) / \mu = t_{16}, \\
 t_{12} &= -2pq^p / \mu, \\
 t_{13} &= -(p^2 - q^p q^s) / \mu = -t_{14}, \\
 t_{15} &= -2pq^s / \mu, \\
 t_{21} &= iq^s / \beta^2 = -t_{23} = -t_{24} = -t_{26}, \\
 t_{31} &= -ip(\Gamma + 2q^p q^s) = t_{36} = t_{41} = t_{46}, \\
 t_{32} &= -4ip^2 q^p, \\
 t_{33} &= -ip(\Gamma - 2q^p q^s) = t_{43} = -t_{34} = -t_{44}, \\
 t_{35} &= -2i\Gamma q^s, \\
 t_{42} &= -2i\Gamma q^p, \\
 t_{45} &= -4ip^2 q^s, \\
 t_{51} &= -iq^p / \beta^2 = t_{53} = t_{54} = -t_{56}, \\
 t_{61} &= -\mu(\Gamma^2 + 4p^2 q^p q^s) = t_{66}, \\
 t_{62} &= -4\mu\Gamma pq^p, \\
 t_{63} &= -\mu(\Gamma^2 - 4p^2 q^p q^s) = -t_{64}, \\
 t_{65} &= -4\mu\Gamma pq^s, \text{ and} \\
 t_{22} &= t_{25} = t_{55} = t_{52} = 0,
 \end{aligned}$$

where  $\Gamma = 2p^2 - 1/\beta^2$ ,  $\mu = \rho\beta^2$ ,  $q^p = (\alpha^{-2} - p^2)^{1/2}$ ,  $q^s = (\beta^{-2} - p^2)^{1/2}$ , and  $p = \sin \theta_p / \alpha = \sin \theta_s / \beta$ .

$T_n^{-1}$  is the inverse matrix of  $T_n$ . The elements of  $T_n^{-1}$  are simply a rearrangement of the elements of  $T_n$ :

$$T_n^{-1} = \begin{bmatrix} t_{61} & t_{51} & t_{31} & t_{31} & t_{21} & t_{11} \\ -t_{65} & 0 & -t_{45} & -t_{35} & 0 & -t_{15} \\ -t_{63} & -t_{51} & -t_{33} & -t_{33} & t_{21} & -t_{13} \\ t_{63} & -t_{51} & t_{33} & t_{33} & t_{21} & t_{13} \\ -t_{62} & 0 & -t_{42} & -t_{32} & 0 & -t_{12} \\ t_{61} & -t_{51} & t_{31} & t_{31} & -t_{21} & t_{11} \end{bmatrix}. \tag{32}$$

### Appendix 2

#### Partial derivation of $E_n$ and $T_n$

In general, the incident angle of angle gather will be controlled within  $90^\circ$ . Therefore, the vertical slowness and  $E_n$  can be expressed as:

$$q^p = (\alpha^{-2} - p^2)^{1/2} = \frac{\cos \theta_p}{\alpha} \tag{33}$$

$$q^s = (\beta^{-2} - p^2)^{1/2} = \frac{\cos \theta_s}{\beta} \tag{34}$$

$$E_n = \text{diag} \left[ e^{-i\omega d_n (\cos \theta_p / \alpha + \cos \theta_s / \beta)} \quad 1 \quad e^{-i\omega d_n (\cos \theta_p / \alpha - \cos \theta_s / \beta)} \quad e^{i\omega d_n (\cos \theta_p / \alpha - \cos \theta_s / \beta)} \quad 1 \quad e^{i\omega d_n (\cos \theta_p / \alpha + \cos \theta_s / \beta)} \right]. \tag{35}$$

The partial derivative of  $E_n$  with respect to the parameters  $M_n$  can be calculated analytically:

$$\frac{\partial E_n}{\partial M_n} = \begin{cases} \frac{\partial E_n}{\partial \alpha} = i\omega d_n \text{diag} \left[ A_\alpha e^{-i\omega d_n (\cos \theta_p / \alpha + \cos \theta_s / \beta)} \quad 0 \quad B_\alpha e^{-i\omega d_n (\cos \theta_p / \alpha - \cos \theta_s / \beta)} \quad -B_\alpha e^{i\omega d_n (\cos \theta_p / \alpha - \cos \theta_s / \beta)} \quad 0 \quad -A_\alpha e^{i\omega d_n (\cos \theta_p / \alpha + \cos \theta_s / \beta)} \right], \\ \frac{\partial E_n}{\partial \beta} = i\omega d_n \text{diag} \left[ A_\beta e^{-i\omega d_n (\cos \theta_p / \alpha + \cos \theta_s / \beta)} \quad 0 \quad B_\beta e^{-i\omega d_n (\cos \theta_p / \alpha - \cos \theta_s / \beta)} \quad -B_\beta e^{i\omega d_n (\cos \theta_p / \alpha - \cos \theta_s / \beta)} \quad 0 \quad -A_\beta e^{i\omega d_n (\cos \theta_p / \alpha + \cos \theta_s / \beta)} \right], \\ \frac{\partial E_n}{\partial \rho} = 0, \end{cases} \tag{36}$$

where

$$\begin{cases} A_\alpha = \left( \frac{\cos \theta_p}{\alpha^2} + \frac{\cos \theta_s}{\alpha\beta} \right), \\ B_\alpha = \left( \frac{\cos \theta_p}{\alpha^2} - \frac{\cos \theta_s}{\alpha\beta} \right), \\ A_\beta = \left( \frac{\cos \theta_p}{\alpha\beta} + \frac{\cos \theta_s}{\beta^2} \right), \\ B_\beta = \left( \frac{\cos \theta_p}{\alpha\beta} - \frac{\cos \theta_s}{\beta^2} \right), \end{cases} \tag{37}$$

in which  $\alpha$  and  $\beta$  are the P- and S-wave velocities at the  $n$ th layer, respectively.

The partial derivatives of  $\mathbf{F}_{n-1}$  and  $\mathbf{F}_n$  with respect to the parameters  $\mathbf{M}_n$  can also be calculated analytically:

$$\frac{\partial \mathbf{F}_n}{\partial \mathbf{M}_n} = \frac{\partial \mathbf{T}_n^{-1}}{\partial \mathbf{M}_n} \mathbf{T}_{n+1}, \quad (38)$$

$$\frac{\partial \mathbf{F}_{n-1}}{\partial \mathbf{M}_n} = \mathbf{T}_{n-1}^{-1} \frac{\partial \mathbf{T}_n}{\partial \mathbf{M}_n}, \quad (39)$$

where

$$\frac{\partial \mathbf{T}_n}{\partial \mathbf{M}_n} = \begin{cases} \frac{\partial \mathbf{T}_n}{\partial \alpha} \\ \frac{\partial \mathbf{T}_n}{\partial \beta} \\ \frac{\partial \mathbf{T}_n}{\partial \rho} \end{cases}. \quad (40)$$

The matrices  $\frac{\partial \mathbf{T}_n}{\partial \alpha}$ ,  $\frac{\partial \mathbf{T}_{n+1}}{\partial \beta}$ , and  $\frac{\partial \mathbf{T}_{n+1}}{\partial \rho}$  contain 16 independent elements.

$$\frac{\partial t_{11}}{\partial \alpha} = -\frac{4}{\alpha} t_{11} = \frac{\partial t_{16}}{\partial \alpha},$$

$$\frac{\partial t_{12}}{\partial \alpha} = -\frac{4}{\alpha} t_{12},$$

$$\frac{\partial t_{13}}{\partial \alpha} = -\frac{4}{\alpha} t_{13} = -\frac{\partial t_{14}}{\partial \alpha},$$

$$\frac{\partial t_{15}}{\partial \alpha} = -\frac{4}{\alpha} t_{15},$$

$$\frac{\partial t_{21}}{\partial \alpha} = -\frac{3}{\alpha} t_{21} = -\frac{\partial t_{23}}{\partial \alpha} = -\frac{\partial t_{24}}{\partial \alpha} = -\frac{\partial t_{26}}{\partial \alpha},$$

$$\frac{\partial t_{31}}{\partial \alpha} = -\frac{3}{\alpha} t_{31} = \frac{\partial t_{36}}{\partial \alpha} = \frac{\partial t_{41}}{\partial \alpha} = \frac{\partial t_{46}}{\partial \alpha},$$

$$\frac{\partial t_{32}}{\partial \alpha} = -\frac{3}{\alpha} t_{32},$$

$$\frac{\partial t_{33}}{\partial \alpha} = -\frac{3}{\alpha} t_{33} = \frac{\partial t_{43}}{\partial \alpha} = -\frac{\partial t_{34}}{\partial \alpha} = -\frac{\partial t_{44}}{\partial \alpha},$$

$$\frac{\partial t_{35}}{\partial \alpha} = -\frac{3}{\alpha} t_{35},$$

$$\frac{\partial t_{42}}{\partial \alpha} = -\frac{3}{\alpha} t_{42},$$

$$\frac{\partial t_{45}}{\partial \alpha} = -\frac{3}{\alpha} t_{45},$$

$$\frac{\partial t_{51}}{\partial \alpha} = -\frac{3}{\alpha} t_{51} = \frac{\partial t_{53}}{\partial \alpha} = \frac{\partial t_{54}}{\partial \alpha} = -\frac{\partial t_{56}}{\partial \alpha},$$

$$\frac{\partial t_{61}}{\partial \alpha} = -\frac{2}{\alpha} t_{61} = \frac{\partial t_{66}}{\partial \alpha},$$

$$\frac{\partial t_{62}}{\partial \alpha} = -\frac{2}{\alpha} t_{62},$$

$$\frac{\partial t_{63}}{\partial \alpha} = -\frac{2}{\alpha} t_{63} = -\frac{\partial t_{64}}{\partial \alpha},$$

$$\frac{\partial t_{65}}{\partial \alpha} = -\frac{2}{\alpha} t_{65},$$

$$\frac{\partial t_{22}}{\partial \alpha} = \frac{\partial t_{25}}{\partial \alpha} = \frac{\partial t_{55}}{\partial \alpha} = \frac{\partial t_{52}}{\partial \alpha} = 0.$$

$$\frac{\partial t_{11}}{\partial \beta} = -\frac{4}{\beta} t_{11} = \frac{\partial t_{16}}{\partial \beta},$$

$$\frac{\partial t_{12}}{\partial \beta} = -\frac{4}{\beta} t_{12},$$

$$\frac{\partial t_{13}}{\partial \beta} = -\frac{4}{\beta} t_{13} = -\frac{\partial t_{14}}{\partial \beta},$$

$$\frac{\partial t_{15}}{\partial \beta} = -\frac{4}{\beta} t_{15},$$

$$\frac{\partial t_{21}}{\partial \beta} = -\frac{3}{\beta} t_{21} = -\frac{\partial t_{23}}{\partial \beta} = -\frac{\partial t_{24}}{\partial \beta} = -\frac{\partial t_{26}}{\partial \beta},$$

$$\frac{\partial t_{31}}{\partial \beta} = -\frac{3}{\beta} t_{31} = \frac{\partial t_{36}}{\partial \beta} = \frac{\partial t_{41}}{\partial \beta} = \frac{\partial t_{46}}{\partial \beta},$$

$$\frac{\partial t_{33}}{\partial \beta} = -\frac{3}{\beta} t_{33} = \frac{\partial t_{43}}{\partial \beta} = -\frac{\partial t_{34}}{\partial \beta} = -\frac{\partial t_{44}}{\partial \beta},$$

$$\frac{\partial t_{35}}{\partial \beta} = -\frac{3}{\beta} t_{35},$$

$$\frac{\partial t_{42}}{\partial \beta} = -\frac{3}{\beta} t_{42},$$

$$\frac{\partial t_{45}}{\partial \beta} = -\frac{3}{\beta} t_{45},$$

$$\frac{\partial t_{51}}{\partial \beta} = -\frac{3}{\beta} t_{51} = \frac{\partial t_{53}}{\partial \beta} = \frac{\partial t_{54}}{\partial \beta} = -\frac{\partial t_{56}}{\partial \beta},$$

$$\frac{\partial t_{61}}{\partial \beta} = -\frac{2}{\beta} t_{61} = \frac{\partial t_{66}}{\partial \beta},$$

$$\frac{\partial t_{62}}{\partial \beta} = -\frac{2}{\beta} t_{62},$$

$$\frac{\partial t_{63}}{\partial \beta} = -\frac{2}{\beta} t_{63} = -\frac{\partial t_{64}}{\partial \beta},$$

$$\frac{\partial t_{65}}{\partial \beta} = -\frac{2}{\beta} t_{65},$$

$$\frac{\partial t_{22}}{\partial \beta} = \frac{\partial t_{25}}{\partial \beta} = \frac{\partial t_{55}}{\partial \beta} = \frac{\partial t_{52}}{\partial \beta} = 0,$$

$$\begin{aligned} \frac{\partial t_{11}}{\partial \rho} &= -\frac{1}{\rho} t_{11} = \frac{\partial t_{16}}{\partial \rho}, \\ \frac{\partial t_{12}}{\partial \rho} &= -\frac{1}{\rho} t_{12}, \\ \frac{\partial t_{13}}{\partial \rho} &= -\frac{1}{\rho} t_{13} = -\frac{\partial t_{14}}{\partial \rho}, \\ \frac{\partial t_{15}}{\partial \rho} &= -\frac{1}{\rho} t_{15}, \\ \frac{\partial t_{21}}{\partial \rho} &= \frac{\partial t_{22}}{\partial \rho} = \frac{\partial t_{23}}{\partial \rho} = \frac{\partial t_{24}}{\partial \rho} = \frac{\partial t_{25}}{\partial \rho} = \frac{\partial t_{26}}{\partial \rho} = 0, \quad \frac{\partial t_{31}}{\partial \rho} = \frac{\partial t_{32}}{\partial \rho} = \frac{\partial t_{33}}{\partial \rho} = \frac{\partial t_{34}}{\partial \rho} = \frac{\partial t_{35}}{\partial \rho} = \frac{\partial t_{36}}{\partial \rho} = 0, \\ \frac{\partial t_{41}}{\partial \rho} &= \frac{\partial t_{42}}{\partial \rho} = \frac{\partial t_{43}}{\partial \rho} = \frac{\partial t_{44}}{\partial \rho} = \frac{\partial t_{45}}{\partial \rho} = \frac{\partial t_{46}}{\partial \rho} = 0, \quad \frac{\partial t_{51}}{\partial \rho} = \frac{\partial t_{52}}{\partial \rho} = \frac{\partial t_{53}}{\partial \rho} = \frac{\partial t_{54}}{\partial \rho} = \frac{\partial t_{55}}{\partial \rho} = \frac{\partial t_{56}}{\partial \rho} = 0, \\ \frac{\partial t_{61}}{\partial \rho} &= \frac{1}{\rho} t_{61} = \frac{\partial t_{66}}{\partial \rho}, \quad \frac{\partial t_{62}}{\partial \rho} = \frac{1}{\rho} t_{62}, \quad \frac{\partial t_{63}}{\partial \rho} = \frac{1}{\rho} t_{63} = -\frac{\partial t_{64}}{\partial \rho}, \quad \frac{\partial t_{65}}{\partial \rho} = \frac{1}{\rho} t_{63}. \end{aligned}$$

## References

- Aki K, Richards PG (1980) Quantitative seismology: theory and methods. W. H. Freeman and Co., San Francisco
- Chen SS, Donoho DL, Saunders MA (2001) Atomic decomposition by basis pursuit. *SIAM Review* 43(1):129–159
- Deng F, McMechan GA (2007) True-amplitude prestack depth migration. *Geophysics* 72(3):S155–S166. <https://doi.org/10.1190/1.2714334>
- Fatti JL, Smith GC, Vail PJ et al (1994) Detection of gas in sandstone reservoirs using AVO analysis: a 3-D seismic case history using the Geostack technique. *Geophysics* 59:1362–1376. <https://doi.org/10.1190/1.1443695>
- Fryer GJ (1980) A slowness approach to the reflectivity method of seismogram synthesis. *Geophys J Int* 63(3):747–758. <https://doi.org/10.1111/j.1365-246X.1980.tb02649.x>
- Fuchs K (1968) The reflection of spherical waves from transition zones with arbitrary depth-dependent elastic moduli and density. *J Phys Earth* 16:27–41. [https://doi.org/10.4294/jpe1952.16.Special\\_27](https://doi.org/10.4294/jpe1952.16.Special_27)
- Fuchs K, Müller G (1971) Computation of synthetic seismograms with the reflectivity method and comparison with observations. *Geophys J Int* 23(4):417–433. <https://doi.org/10.1111/j.1365-246X.1971.tb01834.x>
- Jin S (1999) Characterizing reservoir by using jointly P- and S- wave AVO analyses. *SEG Tech Program Expand Abstr*. <https://doi.org/10.1190/1.1821117>
- Kennett BLN (1974) Reflections, rays, and reverberations. *Bull Seismol Soc Am* 64(6):1685–1696
- Kennett BLN (1983) *Seismic wave propagation in stratified media*. Cambridge University Press, Cambridge. <https://doi.org/10.22459/SWPSM.05.2009>
- Kennett BLN (2009) *Seismic wave propagation in stratified media*. ANU E Press, Canberra
- Larsen JA (1999) AVO inversion by simultaneous P–P and P–S Inversion. M.Sc. thesis, University of Calgary
- Levenberg K (1944) A method for the solution of certain non-linear problems in least squares. *Quart Appl Math* 2:164–168. <https://doi.org/10.1090/qam/1944-02-02>
- Liu HX, Li JY, Chen XH et al (2016) Amplitude variation with offset inversion using the reflectivity method. *Geophysics* 81(4):R185–R195. <https://doi.org/10.1190/geo2015-0332.1>
- Lu J, Yang Z, Wang Y et al (2015) Joint PP and PS AVA seismic inversion using exact Zoeppritz equations. *Geophysics* 80(5):R239–R250. <https://doi.org/10.1190/geo2014-0490.1>
- Lu J, Meng X, Wang Y et al (2016) Prediction of coal seam details and mining safety using multicomponent seismic data: a case history from China. *Geophysics* 81:B149–B165. <https://doi.org/10.1190/geo2016-0009.1>
- Lu J, Wang Y, Chen J et al (2017) Joint anisotropic AVO inversion of PP and PS seismic data. *Geophysics* 83(2):1–83. <https://doi.org/10.1190/geo2016-0516.1>
- Liu W, Wang Y-C, Li J-Y et al (2018) Prestack AVA joint inversion of PP and PS waves using the vectorized reflectivity method. *Appl Geophys* 15(3–4):448–465. <https://doi.org/10.1007/s11770-018-0695-4>
- Luo C, Li XY, Huang GT (2018) Hydrocarbon identification by application of improved sparse constrained inverse spectral decomposition to frequency-dependent AVO inversion. *J Geophys Eng* 15(5):1446–1459. <https://doi.org/10.1088/1742-2140/aab1d6>
- Mahmoudian F, Margrave GF (2004) Three parameter AVO inversion with PP and PS data using offset binning. *SEG Tech Program Expand Abstr*. <https://doi.org/10.1190/1.1851239>
- Mallick S, Frazer LN (1987) Practical aspects of reflectivity modeling. *Geophysics* 52(10):1355–1364. <https://doi.org/10.1190/1.1442248>
- Marquardt DW (1963) An algorithm for least-squares estimation of nonlinear inequalities. *J Soc Ind Appl Math* 11(2):431–441. <https://doi.org/10.1137/0111030>
- Paige CC, Saunders MA (1982) LSQR: an algorithm for sparse linear equations and sparse least squares. *ACM Trans Math Softw (TOMS)* 8(1):43–71. <https://doi.org/10.1145/355984.355989>
- Phinney RA, Odom RI, Fryer GJ (1987) Rapid generation of synthetic seismograms in layered media by vectorization of the algorithm. *Bull Seismol Soc Am* 77(6):2218–2226
- Sen MK, Roy IG (2003) Computation of differential seismograms and iteration adaptive regularization in prestack waveform inversion. *Geophysics* 68(6):2026–2039. <https://doi.org/10.1190/1.1635056>
- Sheen D-H, Tuncay K, Baag C-E et al (2006) Time domain Gauss–Newton seismic waveform inversion in elastic media. *Geophys J Int* 167(3):1373–1384. <https://doi.org/10.1111/j.1365-246X.2006.03162.x>
- Stewart RR (1990) Joint P and P–SV inversion. The CREWES Project research report 2.

- Yuan SY, Liu Y, Zhang Z, Luo CM (2019) Prestack stochastic frequency-dependent velocity inversion with rock-physics constraints and statistical associated hydrocarbon attributes. *IEEE Geosci Remote Sens Lett* 16(1):140–144
- Zhang Y, Zhang G, Bleistein N (2003) True amplitude wave equation migration arising from true amplitude one-way wave equations. *Inverse Probl* 19(5):1113–1138. <https://doi.org/10.1088/0266-5611/19/5/307>





# Three-dimensional angle-domain double-square-root migration in VTI media for the large-scale wide-azimuth seismic data

Chengliang Wu<sup>1</sup> · Bo Feng<sup>1</sup> · Huazhong Wang<sup>1</sup> · Tianzhen Wang<sup>1</sup>

Received: 21 January 2020 / Accepted: 25 May 2020 / Published online: 1 June 2020  
© Institute of Geophysics, Polish Academy of Sciences & Polish Academy of Sciences 2020

## Abstract

With the development of oil and gas exploration, the conventional seismic migration imaging technology based on the isotropic assumption no longer meets our current requirements for high-resolution images. Migration in anisotropic media has become an essential requirement for oil and gas exploration. Marine seismic exploration has gradually entered the wide-azimuth and high-density seismic data acquisition stage. However, even for current large high-performance computer clusters, it is still very difficult to implement pre-stack depth migration based on shot gathers. Thus, we present a double-square-root (DSR) equation based on three-dimensional (3D) pre-stack depth migration in midpoint-offset domain for a wide-azimuth dataset in transversely isotropic media with a vertical symmetry axis (VTI media). Considering VTI media, the DSR migration requires extensive memory and computation; we adopted the phase-shift plus interpolation approach to improve the computational efficiency. Then, we extract the angle-domain common-image gathers (ADCIGs) during DSR migration. For real large-scale seismic data, we designed an effective parallel implementation of 3D DSR migration with ADCIGs outputs. Finally, we applied the proposed angle-domain VTI DSR migration on wide-azimuth SEG/EAGE salt dome-based data and real data from the China South Sea. Numerical and practical data illustrate the effectiveness of the proposed method.

**Keywords** Angle domain · Double-square-root migration · VTI media · Wide-azimuth · Phase-shift plus interpolation

## Introduction

As the target of seismic exploration is turning to complex structures and reservoirs, the broadband, wide-azimuth, and high-density (BWH) seismic data acquisition technology with the long offsets and large observation azimuthal angles is the basic requirement of the high-precision seismic imaging. Developing an accurate seismic imaging technology for wide-azimuth seismic data is a very important issue

(Michell et al. 2006; Barley and Summers 2007; VerWest and Lin 2007; Bouska 2008; Yuan et al. 2019).

The anisotropy phenomenon in subsurface media is ubiquitous. The seismic imaging technology in isotropic (ISO) media can no longer satisfy the high-accuracy reservoir description requirements. For example, ignoring anisotropy-induced distortions due to the difference between the vertical and stacking velocities causes imaging depth errors and ignoring the angle dependence of velocity creates serious problems in imaging dipping reflectors (Tsvankin 2012). Seismic migration imaging in anisotropic media has become an essential requirement for oil and gas exploration. Moreover, when dealing with the wide-azimuth and long-offset seismic data, the anisotropy problem has a more serious impact on migration imaging (Alkhalifah and Larner 1994; Herman and Larner 1995; Yan et al. 2004; Tsvankin 2012; Liu et al. 2014, 2015; Oh and Alkhalifah 2018).

In offshore seismic exploration, no matter the kind of towed-streamer acquisition technology or the kind of wide-azimuth ocean-bottom acquisition technology used, numerous shot gathers are recorded. Even with the current large-scale, high-performance computers, the shot-index migration

---

✉ Bo Feng  
ancd111@163.com

Chengliang Wu  
wuchengliang1990@163.com

Huazhong Wang  
herbhuak@vip.163.com

Tianzhen Wang  
wangtianzhen2011@126.com

<sup>1</sup> Wave Phenomena and Intelligent Inversion Imaging Group (WPI), School of Ocean and Earth Science, Tongji University, Shanghai 200092, China

imaging technology is extremely expensive and unbearable because the computational domain needs to expand with a huge number of zero traces so that the expected reflectors can be imaged in every shot migration. Migration in midpoint-offset domain, which combines all data into one wave-extrapolation procedure, can be an efficient solution (Biondi and Palacharla 1996; Biondi 2002; Alkhalifah et al. 2015).

Angle-domain common-image gathers (ADCIGs) are important outputs for wide-azimuth exploration. The ADCIGs can update the velocity model (Liu and Bleistein 1995; Biondi and Symes 2004) and can be used for extracting the information of angle-dependent reflectivity information (Yan and Xie 2012; Sava et al. 2001; Li et al. 2018). Currently, the image gathers are necessary outputs for any migration algorithms. Therefore, an efficient and robust migration algorithm with the output of image gathers to adapt the anisotropic media and wide-azimuth marine data acquisition is necessary.

Compared with Kirchhoff migration (Gray and May 1994), the wave equation pre-stack depth migration (PSDM) uses accurate wave equations (mainly including one-way and two-way waves), which can obtain high-precision imaging results and is suitable for the current seismic exploration realities (Mulder and Plessix 2003). However, due to the huge computational complexity of reverse-time migration (RTM) (Baysal et al. 1983; Whitmore 1983), the current application of RTM is limited to some local exploration areas for high-precision imaging. The large-scale applications in real industrial exploration are still difficult, especially for TB-level wide-azimuth seismic data. The computational cost is prohibitive in RTM.

Considering computational complexity and imaging accuracy, the one-way wave migration method is a suitable choice. High-quality imaging results can be obtained by using accurate one-way wave equation migration in media with strong heterogeneity, complex structures and dramatic lateral changes velocity. Wavefield continuation migration is an accurate, robust algorithm, even with complex velocity models. There are two main methods for wave equation migration. One is single-square-root (SSR) migration (Reshef 1991; Ke et al. 2004), which continues the up-going and down-going wavefields, respectively, based on the one-way equation and extracts the imaging value by cross-correlating these two wavefields. The second is double-square-root (DSR) migration based on the “survey sinking” concept (Claerbout 1985; Popovici 1996; Biondi and Palacharla 1996; Alkhalifah 2000b; Bevc et al. 2003; de Hoop et al. 2003; Sun et al. 2005; Cheng et al. 2008; Song and Fomel 2011; Alkhalifah et al. 2015).

In the SSR PSDM, because the reflected wave may come from outside the range of offset coverage for both source and receiver wavefields, the extrapolated areas need to be expanded, which increases the computational workload. In addition, the source wavelet is a necessary input in the SSR migration. However, in real data processing, it is impossible to

obtain the complete source wavelet. During the DSR PSDM, only the up-going wavefield is continued. In other words, the shot and receiver locations are sinking at the same time. When the two points coincide, the wavefield value at zero time is regarded as the imaging value. DSR migration does not suffer the problem of migration aperture and has less acquisition footprints. Moreover, it is convenient to produce azimuth-opening angle gathers. Therefore, the DSR migration is more efficient and preferable for real wide-azimuth dataset.

The full 3D DSR PSDM in the midpoint-offset domain is carried out in 5D computational space. In each extrapolation step, it involves totally 3D pre-stack seismic data, which is computationally intensive and difficult to manage. In addition, limited to the narrow azimuth field observation, the offset in the crossline direction of the previous acquired seismic data is small, and the sampling is very sparse, so that the calculation in the crossline direction will bring non-negligible errors. Therefore, in the past decades, due to the limitations of low computer operation speed and immature 3D seismic acquisition techniques, the DSR migration is superseded by the azimuth-moveout (AMO) and common-azimuth migration method (Biondi and Chemingui 1994; Biondi and Palacharla 1996; Biondi et al. 1998; Alkhalifah 2004; Alkhalifah and Biondi 2004). Jin and Wu (1999) and Jin et al. (2002) apply the generalized screen propagators (GSP) (Wu 1994, 1996; de Hoop et al. 2000) to the common-offset depth migration with the DSR operator and use limited azimuthal range of the 3D marine data to reduce the dimensionality. Cheng et al. (2003) proposed a crossline common-offset migration approach for narrow-azimuth datasets. Cheng et al. (2005) further implemented common-azimuth migration utilizing the limited range of the input data volume and the propagation direction of extended wavefields to migrate narrow-azimuth seismic data. However, these methods cannot accurately describe the propagation of seismic waves in the 3D situation and can only be applied to narrow-azimuth dataset. Therefore, they are not suitable approaches for current wide-azimuth datasets. In recent years, high-performance computing and seismic acquisition techniques have both evolved, so now the computational resources are available and the DSR migration has become feasible.

In addition, DSR equation can be implemented in time rather than depth domain (Biondi 2002). However, DSR equation as a time extrapolation operator has an inherent singularity for horizontally traveling waves (Biondi 2002; Duchkov and de Hoop 2009). This singularity can be avoided by using perturbation theory and Shanks transform (Alkhalifah 2012), Padé expansions (Alkhalifah 2013), or limiting the range of wavenumbers treated in a spectral-based extrapolation referred to as the low-rank method (Alkhalifah et al. 2015). In isotropic media, the traveltimes in the midpoint-offset domain can be given by the simple analytical DSR equation (Yilmaz and Claerbout 1980). However, the analytical offset-midpoint travelttime equation for transversely isotropic media does not exist. With stationary

phase approximations and perturbation theory, the analytic approximation representations of the offset-midpoint traveltime in the VTI media can be obtained (Alkhalifah 2000c). Furthermore, the offset-midpoint traveltime is approximately estimated in transversely isotropic media with a horizontal symmetry axis (Hao et al. 2015) and in homogeneous orthorhombic media (Hao et al. 2016).

In this paper, we present a pre-stack depth migration method based on the DSR equation in the midpoint-offset domain in VTI media. In order to reduce the memory and computation, we adopt the phase-shift plus interpolation (PSPI) method to approximate the VTI extrapolation operator. Then, we will discuss how to extract the angle-domain common-image gathers (ADCIGs) during DSR migration. For real large-scale seismic data, we designed a parallel implementation scheme of 3D DSR migration and ADCIGs outputs. Finally, numerical and practical data illustrate the effectiveness of the proposed method.

## Method

In this section, first, we introduce the full 3D DSR wavefield extrapolation operator in VTI media and propose a PSPI method to improve the computational efficiency of wavefield extrapolation operator. Then, we will explain how to extract the ADCIGs during DSR migration. Last, we propose an effective implementation of 3D angle-domain DSR migration for large-scale seismic data.

### Full 3D VTI double-square-root wavefield extrapolation

The qP wave equation, derived by Alkhalifah (2000a), using an acoustic media assumption for P waves in VTI media, yields good kinematics approximation to the familiar elastic wave equation for VTI media. The wavefield solutions obtained using this VTI acoustic wave equation are free of shear waves, which significantly reduces the computation time compared to the elastic wavefield solutions for exploding-reflector-type applications. The 3D qP wave equation is as follows.

$$\begin{aligned} & \frac{\partial^4 p(\mathbf{x}, t)}{\partial t^4} - (1 + 2\varepsilon(\mathbf{x}))V_{p0}^2(\mathbf{x}) \left( \frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2} \right) \frac{\partial^2 p(\mathbf{x}, t)}{\partial t^2} \\ & - V_{p0}^2(\mathbf{x}) \frac{\partial^4 p(\mathbf{x}, t)}{\partial z^2 \partial t^2} \\ & + 2(\varepsilon(\mathbf{x}) - \delta(\mathbf{x}))V_{p0}^4(\mathbf{x}) \left( \frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2} \right) \frac{\partial^2 p(\mathbf{x}, t)}{\partial z^2} = 0, \end{aligned} \quad (1)$$

where  $\mathbf{x}=(x, y, z)$  is the spatial coordinate.  $V_{p0}(\mathbf{x})$  is the vertical P wave velocity,  $\varepsilon(\mathbf{x})$  and  $\delta(\mathbf{x})$  are Thomsen's parameters (Thomsen 1986) and  $p(\mathbf{x}, t)$  is the qP wave wavefield.

The corresponding dispersion relation is (Alkhalifah 2015).

$$\begin{aligned} & \omega^4 - (1 + 2\varepsilon)\omega^2 V_{p0}^2 \mathbf{k}_T^2 - \omega^2 V_{p0}^2 k_z^2 \\ & + 2(\varepsilon - \delta)V_{p0}^4 \mathbf{k}_T^2 k_z^2 = 0, \end{aligned} \quad (2)$$

where  $\omega$  is the angular frequency.  $\mathbf{k}_T = (k_x, k_y)$  is the horizontal wavenumber vector, and  $k_z$  is the  $z$  direction vertical wavenumber.

The DSR migration uses the concept of “survey sinking observation” and takes the wavefield observed on the surface as the boundary condition and the shot and receiver locations are sinking at the same time. During the DSR migration, only the up-going wavefield is continued. According to the form of the one-way wavefield extrapolation operator, the DSR wavefield extrapolation operator in the midpoint-offset domain can be written as follows.

$$p(\mathbf{k}_m, \mathbf{k}_h, z_{j+1}, \omega) = p(\mathbf{k}_m, \mathbf{k}_h, z_j, \omega) e^{ik_z \Delta z}, \quad (3)$$

where  $\mathbf{k}_m = (k_{mx}, k_{my})$  and  $\mathbf{k}_h = (k_{hx}, k_{hy})$  are the wavenumbers of midpoint and half offset, respectively.  $p(\mathbf{k}_m, \mathbf{k}_h, z_{j+1}, \omega)$  and  $p(\mathbf{k}_m, \mathbf{k}_h, z_j, \omega)$  are the extrapolated wavefields at depth  $z_{j+1}$  and  $z_j$ , respectively.  $e^{ik_z \Delta z}$  is the extrapolation propagator.

Based on the Born approximation, the VTI medium is divided into a homogeneous background medium and a perturbed medium.

$$\begin{aligned} s(\mathbf{m}, \mathbf{h}, z) &= s_0(\mathbf{m}, \mathbf{h}, z) + \Delta s(\mathbf{m}, \mathbf{h}, z) \\ \varepsilon(\mathbf{m}, \mathbf{h}, z) &= \varepsilon_0(\mathbf{m}, \mathbf{h}, z) + \Delta \varepsilon(\mathbf{m}, \mathbf{h}, z) \\ \delta(\mathbf{m}, \mathbf{h}, z) &= \delta_0(\mathbf{m}, \mathbf{h}, z) + \Delta \delta(\mathbf{m}, \mathbf{h}, z), \end{aligned} \quad (4)$$

where  $s_0 = 1/V_{p0}$  is the background slowness and  $\Delta s$  is the perturbed slowness.  $\varepsilon_0$  and  $\delta_0$  are the background Thomsen's parameters, and  $\Delta \varepsilon$  and  $\Delta \delta$  are the perturbed Thomsen's parameters, respectively.

The corresponding one-way vertical wavenumber is decomposed into the homogeneous background wavenumber and the perturbed wavenumber. The generalized screen approximation (Wu 1994, 1996; de Hoop et al. 2000) is used to obtain the vertical wavenumber.

$$\tilde{k}_z(\mathbf{m}, \mathbf{h}, z, \mathbf{k}_m, \mathbf{k}_h) = k_{z0}(z, \mathbf{k}_m, \mathbf{k}_h) + k_s(\mathbf{m}, \mathbf{h}, z, \mathbf{k}_m, \mathbf{k}_h), \quad (5)$$

where  $k_{z0}(z, \mathbf{k}_m, \mathbf{k}_h)$  is the vertical wavenumber in the homogeneous background medium and  $k_s(\mathbf{m}, \mathbf{h}, z, \mathbf{k}_m, \mathbf{k}_h)$  is the vertical wavenumber in the perturbed medium. The vertical wavenumber  $k_{z0}(z, \mathbf{k}_m, \mathbf{k}_h)$  can be obtained from the dispersion-relation Eq. (2).

$$k_{z0}(z, \mathbf{k}_m, \mathbf{k}_h) = -\frac{\omega}{v_{p0}} \left[ \left( \frac{1 - (1 + 2\varepsilon_0) \frac{v_{p0}^2}{\omega^2} \left( \frac{\mathbf{k}_m + \mathbf{k}_h}{2} \right)^2}{1 - 2(\varepsilon_0 - \delta_0) \frac{v_{p0}^2}{\omega^2} \left( \frac{\mathbf{k}_m + \mathbf{k}_h}{2} \right)^2} \right)^{1/2} + \left( \frac{1 - (1 + 2\varepsilon_0) \frac{v_{p0}^2}{\omega^2} \left( \frac{\mathbf{k}_m - \mathbf{k}_h}{2} \right)^2}{1 - 2(\varepsilon_0 - \delta_0) \frac{v_{p0}^2}{\omega^2} \left( \frac{\mathbf{k}_m - \mathbf{k}_h}{2} \right)^2} \right)^{1/2} \right]. \tag{6}$$

The  $k_s(\mathbf{m}, \mathbf{h}, z, \mathbf{k}_m, \mathbf{k}_h)$  can be obtained by using the Taylor series expansion. When only the first-order perturbation term is obtained, the approximate wavenumber is as follows (Wu et al. 2007).

$$\tilde{k}_s(\mathbf{m}, \mathbf{h}, z, \mathbf{k}_m, \mathbf{k}_h) = a_{s1}(z, \mathbf{k}_m, \mathbf{k}_h) \Delta s(\mathbf{m}, \mathbf{h}, z) + a_{\varepsilon 1}(z, \mathbf{k}_m, \mathbf{k}_h) \Delta \varepsilon(\mathbf{m}, \mathbf{h}, z) + a_{\delta 1}(z, \mathbf{k}_m, \mathbf{k}_h) \Delta \delta(\mathbf{m}, \mathbf{h}, z), \tag{7}$$

$$\begin{cases} a_{s1}(z, \mathbf{k}_m, \mathbf{k}_h) = \frac{k_0}{k_{z0}} \frac{(1-4(\varepsilon_0-\delta_0)\phi^2+2(1+2\varepsilon_0)(\varepsilon_0-\delta_0)\phi^4)\omega}{(1-2(\varepsilon_0-\delta_0)\phi^2)^2} \\ a_{\varepsilon 1}(z, \mathbf{k}_m, \mathbf{k}_h) = -\frac{k_0}{k_{z0}} \frac{(1+2\delta_0)\phi^4}{(1-2(\varepsilon_0-\delta_0)\phi^2)^2} k_0 \\ a_{\delta 1}(z, \mathbf{k}_m, \mathbf{k}_h) = -\frac{k_0}{k_{z0}} \frac{(1-(1+2\varepsilon_0)\phi^2)\phi^2}{(1-2(\varepsilon_0-\delta_0)\phi^2)^2} k_0, \end{cases} \tag{7a}$$

where  $k_0 = \frac{\omega}{v_0} = \omega s_0$  and  $\phi = \frac{k_x}{k_0}$ .  $a_{s1}$ ,  $a_{\varepsilon 1}$  and  $a_{\delta 1}$  are the Taylor series expansion coefficients of the anisotropic parameter perturbations with respect to slowness perturbations under the generalized screen approximation, respectively.

Using the approximation  $e^x \approx 1 + x$ , the final DSR extrapolated wavefield in VTI media can be realized as follows:

$$p(\mathbf{m}, \mathbf{h}, z_{j+1}, \omega) = p_{ps}(\mathbf{m}, \mathbf{h}, z_{j+1}, \omega) + p_{\Delta s}(\mathbf{m}, \mathbf{h}, z_{j+1}, \omega) + p_{\Delta \varepsilon}(\mathbf{m}, \mathbf{h}, z_{j+1}, \omega) + p_{\Delta \delta}(\mathbf{m}, \mathbf{h}, z_{j+1}, \omega), \tag{8}$$

$$p_{ps}(\mathbf{m}, \mathbf{h}, z_{j+1}, \omega) = IF \left\{ e^{ik_{z0}(z_j)\Delta z} F \left\{ e^{i\omega \Delta s \Delta z} p(\mathbf{m}, \mathbf{h}, z_j, \omega) \right\} \right\}, \tag{8a}$$

$$p_{\Delta s}(\mathbf{m}, \mathbf{h}, z_{j+1}, \omega) = IF \left\{ (a_{s1_s} - \omega) e^{ik_{z0}(z_j)\Delta z} F \left\{ e^{i\omega \Delta s \Delta z} i \Delta s_s \Delta z p(\mathbf{m}, \mathbf{h}, z_j, \omega) \right\} \right\} + IF \left\{ (a_{s1_g} - \omega) e^{ik_{z0}(z_j)\Delta z} F \left\{ e^{i\omega \Delta s \Delta z} i \Delta s_g \Delta z p(\mathbf{m}, \mathbf{h}, z_j, \omega) \right\} \right\}, \tag{8b}$$

$$p_{\Delta \varepsilon}(\mathbf{m}, \mathbf{h}, z_{j+1}, \omega) = IF \left\{ a_{\varepsilon 1_s} e^{ik_{z0}(z_j)\Delta z} F \left\{ e^{i\omega \Delta s \Delta z} i \Delta \varepsilon_s \Delta z p(\mathbf{m}, \mathbf{h}, z_j, \omega) \right\} \right\} + IF \left\{ a_{\varepsilon 1_g} e^{ik_{z0}(z_j)\Delta z} F \left\{ e^{i\omega \Delta s \Delta z} i \Delta \varepsilon_g \Delta z p(\mathbf{m}, \mathbf{h}, z_j, \omega) \right\} \right\}, \tag{8c}$$

$$p_{\Delta \delta}(\mathbf{m}, \mathbf{h}, z_{j+1}, \omega) = IF \left\{ a_{\delta 1_s} e^{ik_{z0}(z_j)\Delta z} F \left\{ e^{i\omega \Delta s \Delta z} i \Delta \delta_s \Delta z p(\mathbf{m}, \mathbf{h}, z_j, \omega) \right\} \right\} + IF \left\{ a_{\delta 1_g} e^{ik_{z0}(z_j)\Delta z} F \left\{ e^{i\omega \Delta s \Delta z} i \Delta \delta_g \Delta z p(\mathbf{m}, \mathbf{h}, z_j, \omega) \right\} \right\}, \tag{8d}$$

where F represents the Fourier transform and IF represents the inverse Fourier transform.  $p_{k_{z0}}$  is the extrapolated wavefield in the background medium.  $p_{\Delta s}$ ,  $p_{\Delta \varepsilon}$  and  $p_{\Delta \delta}$  are the extrapolated wavefield in the perturbed medium with respect to the slowness,  $\varepsilon$ , and  $\delta$ , perturbations, respectively.

In the midpoint-offset domain, the source and receiver wavefields are simultaneously extrapolated along the vertical depth axis using the DSR operator. The migrated results are extracted by applying the Claerbout’s (1985) imaging condition.

$$I(\mathbf{m}, z) = I(t = 0, \mathbf{k}_m, \mathbf{h}) = 0, z + \Delta z = \int d\omega \int d\mathbf{k}_h e^{ik_z(\omega, \mathbf{k}_m, \mathbf{k}_h)\Delta z} P(\mathbf{k}_m, \mathbf{k}_h, z, \omega), \tag{9}$$

$$P(\mathbf{k}_m, \mathbf{k}_h, z, \omega) = \int d\mathbf{k}_m e^{-i\omega \mathbf{k}_m} \int d\mathbf{k}_h e^{-i\omega \mathbf{k}_h} p(\mathbf{m}, \mathbf{h}, z, \omega), \tag{9a}$$

where  $I(\mathbf{m}, z)$  is the imaging result.  $P(\mathbf{k}_m, \mathbf{k}_h, z, \omega)$  is the 2D Fourier transform of the extrapolated wavefield  $p(\mathbf{m}, \mathbf{h}, z, \omega)$ . The integrals of  $\omega$  and  $k_h$  reflect the imaging condition:  $t = 0$  and  $h = 0$ .

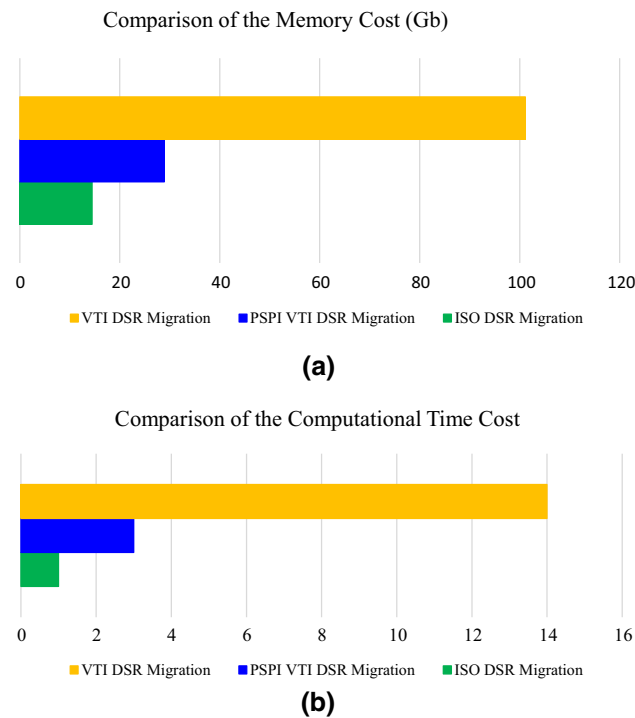
### Full 3D VTI DSR migration based on phase-shift plus interpolation

From Eq. (8) and (9), we see that, in order to implement the full 3D VTI wavefield extrapolation, both 4D FFT and 4D IFFT need to be applied 7 times at each extrapolated depth. Because the extrapolation operator is a hybrid domain operator, seven 4D arrays are required to store the wavefield. Therefore, the computational cost and memory requirements are quite high, especially for real wide-azimuth seismic data.

A small-scale wide-azimuth salt dome model is used to analyze the computational complexity and memory requirements. Table 1 shows the parameters of the salt dome model, and the computational time cost and memory requirements are shown in Fig. 1. As we can see, in the salt dome case, the needed random-access memory (RAM) is less and affordable for the ISO DSR migration. However, the RAM

**Table 1** The parameters for the salt dome model

Parameters	Values	Parameters	Values (m)
Inline range	(750 to 12,750 m)	Dmx	30
Crossline range	(4710 to 8310 m)	Dmy	30
Offset range of inline	(−3000 to 3000 m)	Dhx	30
Offset range of crossline	(−3000 to 3000 m)	Dhy	30



**Fig. 1** Comparison of **a** major random-access memory requirement and **b** computational time cost for the salt dome model

requirement for the VTI DSR migration is several times higher than the ISO DSR migration. The RAM requirement of the VTI DSR migration becomes extremely large, so that it is not feasible in practical application, even with the current high-performance computers. In addition, compared with the ISO DSR migration, the VTI DSR migration has low computational efficiency.

In order to reduce the RAM requirements and improve the computational efficiency, we propose using phase-shift plus interpolation (PSPI) method (Gazdag and Sguazzero 1984) to realize the VTI DSR migration. Using perturbation theory, we know that, compared with the anisotropic parameter perturbation, the velocity perturbation will affect the imaging results more significantly. Therefore, the new VTI extrapolation operator is obtained by considering the influence of the slowness perturbation and neglecting the

influence of the anisotropic parameter perturbation. Equation (5) becomes:

$$\tilde{k}_z(\mathbf{m}, \mathbf{h}, z, \mathbf{k}_m, \mathbf{k}_h) = k_{z0}(z, \mathbf{k}_m, \mathbf{k}_h) + k_{\Delta s}(\mathbf{m}, \mathbf{h}, z, \mathbf{k}_m, \mathbf{k}_h), \tag{10}$$

where the vertical background wavenumber  $k_{z0}(z, \mathbf{k}_m, \mathbf{k}_h)$  is given by the DSR Eq. (6). The first-order expansion formula of the scattering operator  $k_{\Delta s}(\mathbf{m}, \mathbf{h}, z, \mathbf{k}_m, \mathbf{k}_h)$  is obtained as follows.

$$k_{\Delta s}(\mathbf{m}, \mathbf{h}, z, \mathbf{k}_m, \mathbf{k}_h) = a_{s1_s}(z, \mathbf{k}_m, \mathbf{k}_h) \Delta s_s(\mathbf{m}, \mathbf{h}, z) + a_{s1_g}(z, \mathbf{k}_m, \mathbf{k}_h) \Delta s_g(\mathbf{m}, \mathbf{h}, z), \tag{11}$$

$$a_{s1_s}(z, \mathbf{k}_m, \mathbf{k}_h) = \frac{k_0}{k_{z0}} \frac{(1 - 4(\epsilon_0 - \delta_0)\varphi_s + 2(1 + 2\epsilon_0)(\epsilon_0 - \delta_0)\varphi_s^4)}{(1 - 2(\epsilon_0 - \delta_0)\varphi_s)^2} \omega, \tag{11a}$$

$$a_{s1_g}(z, \mathbf{k}_m, \mathbf{k}_h) = \frac{k_0}{k_{z0}} \frac{(1 - 4(\epsilon_0 - \delta_0)\varphi_g + 2(1 + 2\epsilon_0)(\epsilon_0 - \delta_0)\varphi_g^4)}{(1 - 2(\epsilon_0 - \delta_0)\varphi_g)^2} \omega, \tag{11b}$$

$$\varphi_g = \frac{v_{p0_g}^2}{\omega^2} \left( \left( \frac{k_{mx} + k_{hx}}{2} \right)^2 + \left( \frac{k_{my} + k_{hy}}{2} \right)^2 \right);$$

$$\varphi_s = \frac{v_{p0_s}^2}{\omega^2} \left( \left( \frac{k_{mx} - k_{hx}}{2} \right)^2 + \left( \frac{k_{my} - k_{hy}}{2} \right)^2 \right), \tag{11c}$$

where  $\Delta s_s$  and  $\Delta s_g$  are the slowness perturbations for the source and receiver positions, respectively.

Then, we use the PSPI method to compensate the influence of the anisotropic parameter perturbations and achieve a more accurate migration result. The PSPI method can be implemented in frequency-wavenumber domain. The basic principle of the method is to extrapolate using two or more reference anisotropic parameters downward to obtain multiple reference wavefields. Then, according to the relationship between actual migration parameters and the reference parameters, we obtain the final extrapolated wavefield with the aid of the appropriate interpolation method. The DSR extrapolated wavefield can be written in the following form.

$$p(\mathbf{k}_m, \mathbf{k}_h, z_{j+1}, \omega) = p(\mathbf{k}_m, \mathbf{k}_h, z_j, \omega) e^{i(k_s(\epsilon_s, \delta_s) + k_r(\epsilon_r, \delta_r))dz}, \tag{12}$$

where  $\epsilon_s, \delta_s$  and  $\epsilon_r, \delta_r$  are the anisotropic parameters for the source and receiver wavefields, respectively. However, the computational cost of the interpolation process will become relatively high when simultaneously taking the influence of  $\epsilon$  and  $\delta$  into consideration, according to the fact that the

parameter  $\varepsilon$  is usually associated with the parameter  $\delta$  and the parameter  $\varepsilon$  has a greater impact on imaging than the parameter  $\delta$ . Therefore, in the wavefield interpolation procedure, we just consider the effect of the parameter  $\varepsilon$ . Finally, the wavefield interpolation is implemented in the following manner.

Step 1: We calculate the phase shift by an ISO DSR propagator to obtain the wavefield  $p_z^1(\mathbf{k}_m, \mathbf{k}_h, z_j, \omega)$ ;

Step 2: We use the average parameter of  $(\varepsilon^{avg}, \delta^{avg})$  to get the phase shift result of the wavefield  $p_z^2(\mathbf{k}_m, \mathbf{k}_h, z_j, \omega)$  with the proposed VTI DSR propagator, and obtain the wavefield  $p_z^3(\mathbf{k}_m, \mathbf{k}_h, z_j, \omega)$  using the parameter  $(\varepsilon^{max}, \delta^{max})$ ;

Step 3: Finally, we adopt the following formula for interpolation to get the final DSR wavefield.

$$p(\mathbf{k}_m, \mathbf{k}_h, z_j, \omega) = a_z^1 \times p_z^1(\mathbf{k}_m, \mathbf{k}_h, z_j, \omega) + a_z^2 \times p_z^2(\mathbf{k}_m, \mathbf{k}_h, z_j, \omega) + a_z^3 \times p_z^3(\mathbf{k}_m, \mathbf{k}_h, z_j, \omega). \tag{13}$$

These coefficient parameters  $a_z^1$ ,  $a_z^2$  and  $a_z^3$  can be obtained by the relationship.

$$a_z^1 = \frac{(\varepsilon_z^{sr} - \varepsilon_z^{avg})(\varepsilon_z^{sr} - \varepsilon_z^{max})}{(0 - \varepsilon_z^{avg})(0 - \varepsilon_z^{max})}, \tag{13a}$$

$$a_z^2 = \frac{(\varepsilon^{sr} - 0)(\varepsilon^{sr} - \varepsilon^{max})}{(\varepsilon^{avg} - 0)(\varepsilon^{avg} - \varepsilon^{max})}, \tag{13b}$$

$$a_z^3 = \frac{(\varepsilon^{sr} - 0)(\varepsilon^{sr} - \varepsilon^{avg})}{(\varepsilon^{max} - 0)(\varepsilon^{max} - \varepsilon^{avg})}, \tag{13c}$$

$$\varepsilon^{sr} = \frac{\varepsilon^s + \varepsilon^r}{2}. \tag{13d}$$

In addition, Alkhalifah (1998, 2000a) and Alkhalifah et al. (2001) proposed a new parameterization representation in terms of just two parameters in anisotropic media.

$$v = v_{p0} \sqrt{1 + 2\delta}, \tag{14}$$

$$\eta = \frac{\varepsilon - \delta}{1 + 2\delta}, \tag{15}$$

where  $v$  is the velocity and  $\eta$  is the anisotropy coefficient. The approximation is far more accurate than the weak-anisotropy (Thomsen 1986) or the small-angle approximation (Cohen 1997), while it can simplify the equations (Alkhalifah 1998).

Therefore, under this representation, the vertical background wavenumber  $k_{z0}(z, \mathbf{k}_m, \mathbf{k}_h)$  and the perturbed wavenumber  $k_{\Delta s}(\mathbf{m}, \mathbf{h}, z, \mathbf{k}_m, \mathbf{k}_h)$  become

$$k_{z0}(z, \mathbf{k}_m, \mathbf{k}_h) = -\frac{\omega}{v_{p0}} \left[ \frac{1 - (1 + 2\eta_0) \frac{v_0^2}{\omega^2} \left(\frac{\mathbf{k}_m + \mathbf{k}_h}{2}\right)^2}{1 - 2\eta_0 \frac{v_0^2}{\omega^2} \left(\frac{\mathbf{k}_m + \mathbf{k}_h}{2}\right)^2} \right]^{1/2} + \left[ \frac{1 - (1 + 2\eta_0) \frac{v_0^2}{\omega^2} \left(\frac{\mathbf{k}_m - \mathbf{k}_h}{2}\right)^2}{1 - 2\eta_0 \frac{v_0^2}{\omega^2} \left(\frac{\mathbf{k}_m - \mathbf{k}_h}{2}\right)^2} \right]^{1/2}. \tag{16}$$

$$k_{\Delta s}(\mathbf{m}, \mathbf{h}, z, \mathbf{k}_m, \mathbf{k}_h) = a_{s1-s}(z, \mathbf{k}_m, \mathbf{k}_h) \Delta s_s(\mathbf{m}, \mathbf{h}, z) + a_{s1-g}(z, \mathbf{k}_m, \mathbf{k}_h) \Delta s_g(\mathbf{m}, \mathbf{h}, z), \tag{17}$$

$$a_{s1-s}(z, \mathbf{k}_m, \mathbf{k}_h) = \frac{k_0}{k_{z0}} \frac{\left(1 - 8\eta_0 \frac{v_0^2}{\omega^2} \tilde{\varphi}_s + 4\eta_0(1 + 2\eta_0) \frac{v_0^4}{\omega^4} \frac{v_{p0-s}^4}{\omega^4} \tilde{\varphi}_s^4\right)}{\left(1 - 2\eta_0 \frac{v_0^2}{\omega^2} \tilde{\varphi}_s\right)^2} \omega, \tag{17a}$$

$$a_{s1-g}(z, \mathbf{k}_m, \mathbf{k}_h) = \frac{k_0}{k_{z0}} \frac{\left(1 - 8\eta_0 \frac{v_0^2}{\omega^2} \tilde{\varphi}_g + 4\eta_0(1 + 2\eta_0) \frac{v_0^4}{\omega^4} \frac{v_{p0-s}^4}{\omega^4} \tilde{\varphi}_g^4\right)}{\left(1 - 2\eta_0 \frac{v_0^2}{\omega^2} \tilde{\varphi}_g\right)^2} \omega, \tag{17b}$$

$$\begin{cases} \tilde{\varphi}_g = \left( \left( \frac{k_{mx} + k_{hx}}{2} \right)^2 + \left( \frac{k_{my} + k_{hy}}{2} \right)^2 \right) \\ \tilde{\varphi}_s = \left( \left( \frac{k_{mx} - k_{hx}}{2} \right)^2 + \left( \frac{k_{my} - k_{hy}}{2} \right)^2 \right) \end{cases}, \tag{17c}$$

where  $\Delta s_s$  and  $\Delta s_g$  are the perturbed slowness for the source and receiver positions, respectively.

The DSR extrapolated wavefield with the PSPI method then has the following form.

$$p(\mathbf{k}_m, \mathbf{k}_h, z_j, \omega) = a_z^1 \times p_z^1(\mathbf{k}_m, \mathbf{k}_h, z_j, \omega) + a_z^2 \times p_z^2(\mathbf{k}_m, \mathbf{k}_h, z_j, \omega) + a_z^3 \times p_z^3(\mathbf{k}_m, \mathbf{k}_h, z_j, \omega), \tag{18}$$

$$a_z^1 = \frac{(\eta_z^{sr} - \eta_z^{avg})(\eta_z^{sr} - \eta_z^{max})}{(0 - \eta_z^{avg})(0 - \eta_z^{max})}, \tag{18a}$$

$$a_z^2 = \frac{(\eta^{sr} - 0)(\eta^{sr} - \eta^{max})}{(\eta^{avg} - 0)(\eta^{avg} - \eta^{max})}, \tag{18b}$$

$$a_z^3 = \frac{(\eta^{sr} - 0)(\eta^{sr} - \eta^{avg})}{(\eta^{max} - 0)(\eta^{max} - \eta^{avg})}, \tag{18c}$$

$$\eta^{sr} = \frac{\eta^s + \eta^r}{2}, \tag{18d}$$

where the wavefield  $p_z^1(\mathbf{k}_m, \mathbf{k}_h, z_j, \omega)$  is obtained using the ISO DSR propagator and the wavefields  $p_z^2(\mathbf{k}_m, \mathbf{k}_h, z_j, \omega)$  and

$p_z^3(\mathbf{k}_m, \mathbf{k}_h, z_j, \omega)$  are calculated with the parameters of  $\eta^{\text{avg}}$  and  $\eta^{\text{max}}$ , respectively.

### Extraction of the ADCIGs during 3D DSR migration

The quality and computational efficiency of image gathers are the main objectives of the migration procedure. The ray-based migration methods, such as Kirchhoff and Beam migrations, can conveniently produce ADCIGs, at the cost of a decreasing imaging accuracy (Cai et al. 2013; Liu et al. 2018). On the other hand, RTM can produce relatively accurate ADCIGs, but the computational cost is much higher (Xu et al. 2011; Jin et al. 2014; Hu et al. 2015; Wu et al. 2019). However, the 3D-DSR equation can output ADCIGs which has a good balance between the imaging accuracy and the computational cost (Sava and Fomel 2005; Biondi 2007; Sava and Vlad 2011; Sava and Alkhalifah 2013).

In the midpoint-offset domain, we can obtain the imaging result by integrating the extrapolated wavefield along the dimensions  $\omega$  and  $k_h$  in Eq. (9). If we don't stack the offset information in the extrapolated procedure, the offset-domain common-image gathers (ODCIGs) can be extracted through the imaging condition as follows.

$$I(\mathbf{m}, \mathbf{k}_h, z) = \int d\omega e^{ik_z(\omega, \mathbf{k}_m, \mathbf{k}_h)z} P(\omega, \mathbf{k}_m, \mathbf{k}_h, z). \quad (19)$$

We generate and store the local-offset gathers using the ODCIGs imaging condition. The local-offset gathers can be transformed into the Fourier domain.

$$I(\mathbf{m}, \mathbf{k}_h, z) \rightarrow I(\mathbf{m}, \mathbf{k}_h, \mathbf{k}_z), \quad (20)$$

where  $\rightarrow$  represents the mapping process.

Then the azimuth-opening ADCIGs can be extracted through radial trace construction.

$$I(\mathbf{m}, \mathbf{k}_h, \mathbf{k}_z) \rightarrow I(\theta, \gamma, \mathbf{m}, \mathbf{k}_z). \quad (21)$$

The relationship between the angle and wavenumber is as follows (Sava and Fomel 2003).

$$\tan \theta = -\text{sign}(k_{hx}) \frac{|\mathbf{k}_h|}{k_z}, \quad (22a)$$

$$\tan \gamma = \frac{k_{hx}}{k_{hy}}, \quad (22b)$$

where  $\theta$  is the reflection angle and  $\gamma$  is the azimuthal angle shown in Fig. 2.

The final angle imaging result can be obtained by inverse Fourier transform.

$$I(\theta, \gamma, \mathbf{m}, \mathbf{k}_z) \rightarrow I(\theta, \gamma, \mathbf{m}, z). \quad (23)$$

The above process is called post-migration angle gathers extraction implemented in the imaging space (Fomel 2004; Alkhalifah and Fomel 2011), which does not require a large storage space and extra computation. Therefore, the associated ADCIGs can be applied for subsequent velocity inversion (Biondi and Symes 2004).

In addition, in full 3D DSR migration, when the wavefield is transformed into the frequency-wavenumber domain, we can also conveniently extract the offset ray parameter CIGs using the mapping formula (Sava and Fomel 2003).

$$\mathbf{p}_h = \frac{\mathbf{k}_h}{\omega}, \quad (24)$$

where  $\mathbf{p}_h = (p_{hx}, p_{hy})$  is the offset wavevector. Then, the image gathers are extracted through radial trace transforms in the Fourier domain in the data space. The conversion to angle gathers is implemented at every extrapolation step and is less sensitive to inaccuracies in the location of sharp velocity boundaries.

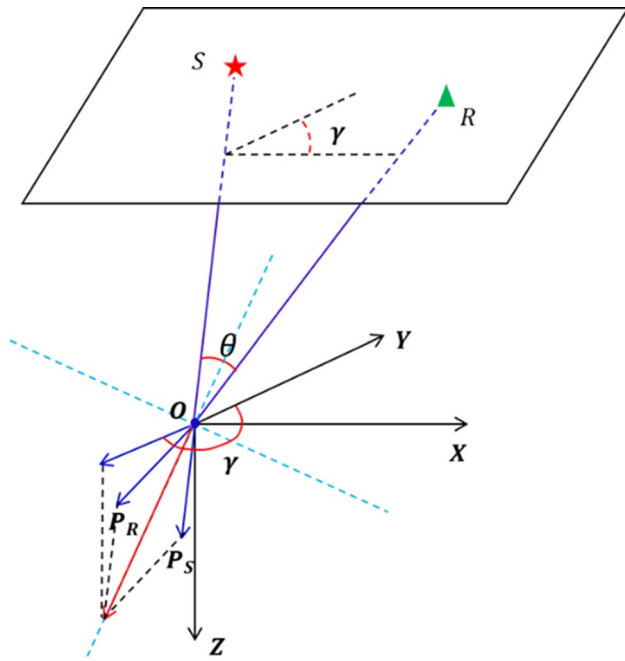
### An effective scheme for 3D angle-domain DSR migration

In order to achieve best computational efficiency with the available computer cluster, the realization of 3D DSR migration and ADCIGs output must adapt to the hardware structure of computer cluster. Multiple nodes with large single-node memory with a global disk and limited local disks, and multiple input/output (I/O) channels are the main characteristics of the cluster. Considering the memory usage, seismic data I/O, imaging results I/O and the migration imaging accuracy, we designed a parallel implementation scheme for 3D DSR migration and ADCIGs outputs considering a large-scale seismic data volume, which is shown in Fig. 3. The scheme mainly includes three parts as follows.

(1) Preprocessing of migration data.

The main purpose of this part is to convert the input data into the frequency domain to prepare them for DSR migration. Generally, the time dimension is the fastest dimension (innermost loop) of the input seismic data. However, the outermost loop in DSR migration imaging is the frequency dimension. Therefore, the input data must be reordered to adapt to the migration algorithm. In order to adapt to the inconsistent observation systems, which the order of data is inconsistent, the following strategy is implemented.

Firstly, we use the header index information to sort the data and store the coordinates of midpoint, offset and storage location. Then, the required data are sorted by a two-step method (see formula 25). Taking 3D CMP seismic data as an example, the input data are five-dimensional, and the



**Fig. 2** Schematic representation of the azimuth angle  $\gamma$  and opening angle  $\theta$  at the common imaging point  $O$  by source and receiver ray-parameter vectors ( $P_S$  and  $P_R$ )

directly transformed procedure is memory-consuming. However, this process can be easily implemented without consuming too much storage, if we adopt the two-step method to realize this process.

$$(t, hx, hy, mx, my) \rightarrow (hx, hy, \omega, mx, my) \rightarrow (hx, hy, mx, my, \omega) \tag{25}$$

In the meanwhile, the transformed data are placed on the right order depending on the recorded information. Furthermore, we use the MPI and parallel I/O to improve efficiency. Moreover, in order to reduce the final storage space, we only store the position where the CMP point is illuminated by the source and receiver. The storage format now becomes as follows.

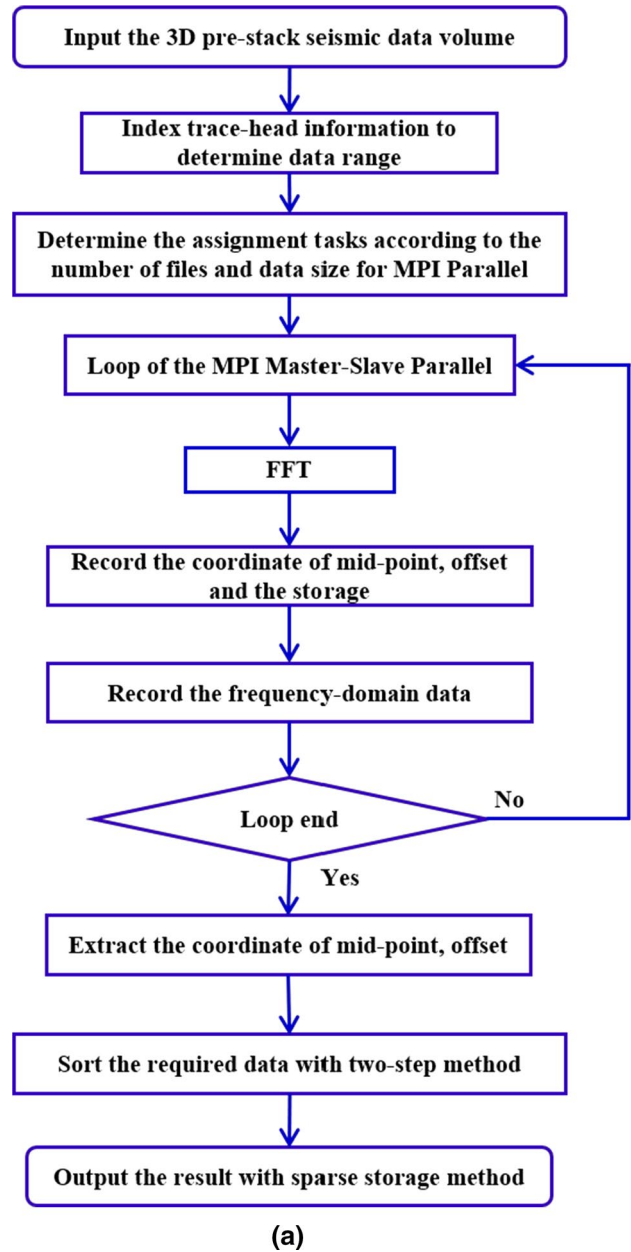
$$(hx, hy, mx, my, \omega) \Rightarrow (nfold\_max, mx, my, \omega) \tag{26}$$

where the *nfold\_max* is the maximum illuminated number. For example, in the wide-azimuth salt dome example, if the frequency range is 2–60 Hz and the time sampling is 4 ms, the origin storage size is nearly 7000 GB. However, the new storage procedure only requires 211.75 GB, because the other approach includes a lot of zero elements.

(2) DSR migration.

The main purpose of this part is to extrapolate the input data in the frequency domain with the anisotropic extrapolation operator and output the imaging result and the subsurface migration gathers. The loop process of DSR migration is shown in Fig. 3b. Firstly, we determine

the output crossline range based on the size of computer memory and apply the MPI master–slave parallel loop for all frequencies and crossline range. The DSR extrapolation operator requires Fourier transform from the midpoint–offset domain to wavenumber domain. The sampling in the wavenumber domain is often irregular and under-sampled. Binning and regularization are implemented using the migration operator and taking into account the image grid spacing to alleviate the aliasing in the wavenumber domain. The OpenMP parallelism is used in the



**Fig. 3** Flowcharts for the 3D angle-domain DSR migration, **a** pre-processing of the migration data, **b** DSR migration, **c** imaging gathers output



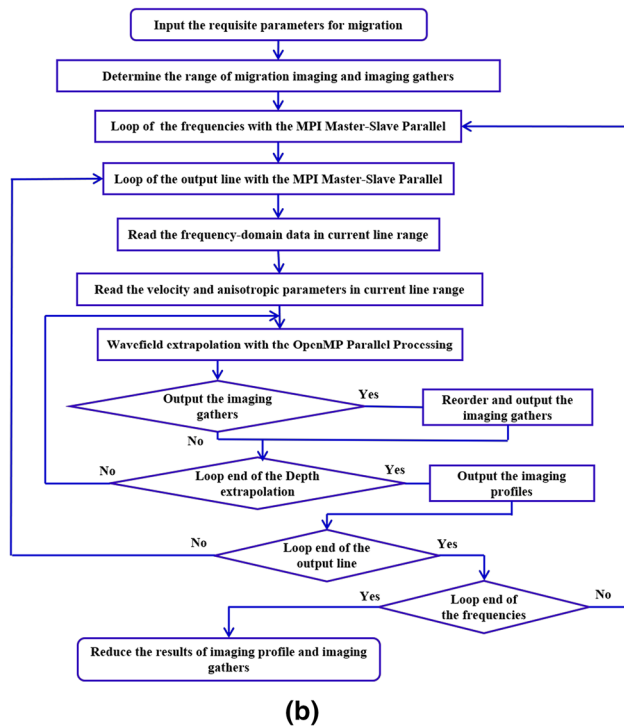


Fig. 3 (continued)

extrapolation step, and the 4D Fourier transform between the space and wavenumber domains is implemented by the multi-threaded FFTW (Frigo and Steven 2005) to improve the computational efficiency. In order to avoid the direct output of 5D image gathers and reduce the required memory, we store the image gathers in the local disk of the current node during each extrapolation step. The output image from a single node is stored on the current local disk after the depth extrapolation. Finally, all images and image gathers are superimposed on the global disk.

### (3) The image gathers output.

Converting the subsurface migration gathers into azimuth-opening ADCIGs is easy to implement and thus allow for an ease in the implementation of the parallel strategy shown in Fig. 3c. The interpolation and regularization are used to improve the quality of angle gathers. Compared with the DSR migration step, the computational time cost and memory requirements in the mapping process are negligible.

## Numerical examples

In this section, we illustrate the effectiveness of the proposed VTI DSR migration method on both synthetic and field data examples.

In the first example, we use the wide-azimuth SEG–EAGE salt dome model to demonstrate the validity of the proposed

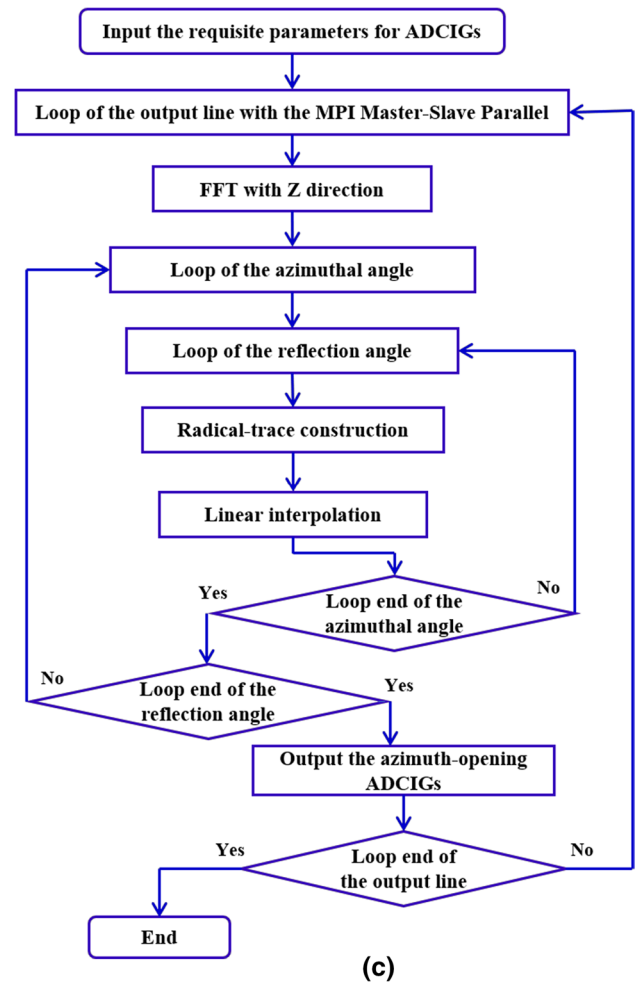
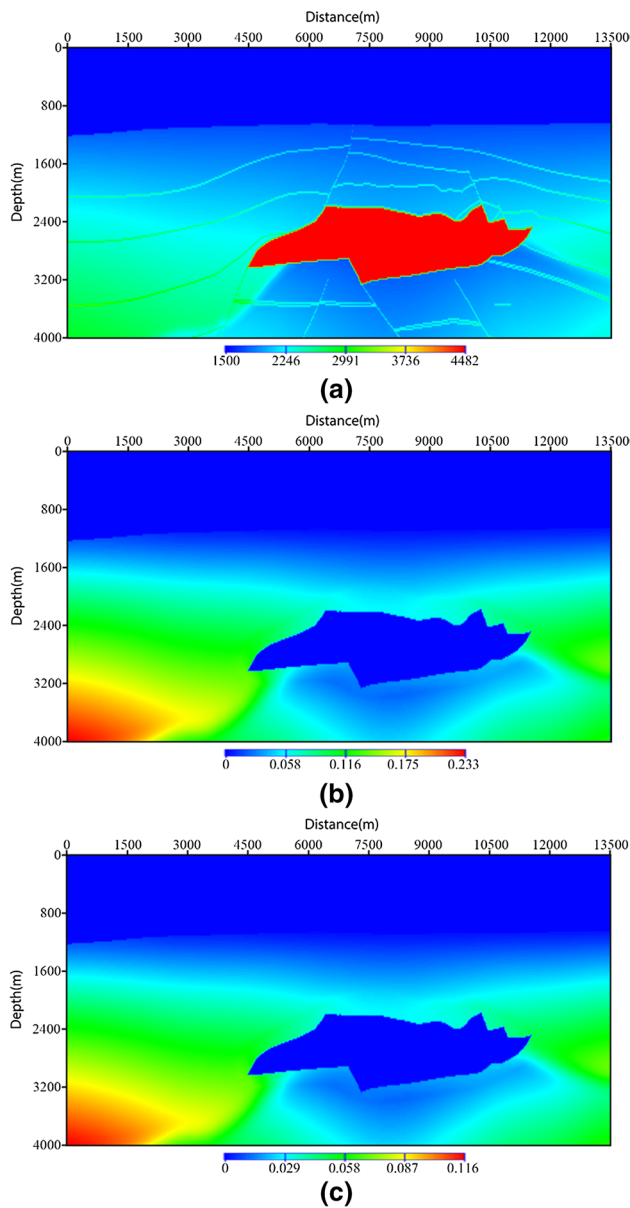


Fig. 3 (continued)

3D VTI DSR migration. The size of the velocity model and the anisotropic parameters are 13.5 km and 13.5 km in  $X$  and  $Y$  directions with a 30-m grid interval and 5 km in the  $Z$  direction with a 20-m grid interval. The parameters of the acquisition system are as follows: The inline range of shots is 750–12,750 m, and the crossline range of shots is 4710–8310 m. The interval of between shots is 120 m in the  $X$  direction and 360 m in the  $Y$  direction. In every shot gathers, the offset range of geophone is  $-3000$  to  $3000$  m with a 30 m grid interval in both directions. The synthetic data are generated with a finite-difference forward modeling method. The source function is a Ricker wavelet, with a dominant frequency of 20 Hz, and the time of the shot gathers extends to 8 s with a 4-ms sampling interval. The frequency range used in migration process is 2–60 Hz. Figure 4a is an inline section of the velocity model, and Fig. 4b–c shows the Thomsen's anisotropic parameters ( $\epsilon$  and  $\delta$ ) along the same inline where  $\epsilon$  ranges between 0 and 0.233, and  $\delta$  spans the range 0–0.116. The  $\epsilon$  and  $\delta$  values are large on both sides of the salt dome; meanwhile, inside salt dome the  $\epsilon$  and  $\delta$  values are small. Ignoring the



**Fig. 4** Inline profile of the SEG–EAGE Salt dome model for **a** migrated velocity model and **b** Thomsen's anisotropic parameters ( $\epsilon$ ), and **c** Thomsen's anisotropic parameters ( $\delta$ )

anisotropic effects, the imaging result with an ISO extrapolation operator is shown in Fig. 5a. Figure 5b shows the migrated result with the proposed VTI wavefield extrapolation operator. As we can see, because the influence of anisotropic parameters is neglected, the result of the isotropic migration is generally poor. However, the VTI imaging result is better focused and clearer than the ISO result.

Moreover, the difference can be better analyzed and more obvious in the image gathers. In this example, we extract the offset ray parameter CIGs to illustrate the effectiveness of the proposed VTI DSR migration. Figure 6a, b shows the

common-image gathers by the ISO migration and the proposed VTI migration, respectively. Because we adopt the correct migration velocity model and anisotropic parameters in the migration process, the image gathers should be even for the reflection events. However, in the ISO imaging result, the reflection events are uneven and there is some curvature such as the red circles on the left side of salt dome. Compensating for the influence of anisotropic parameters in the proposed VTI image gathers produces reflection events that are flat and have better focused. The gathers can provide more useful input information for velocity tomography, which picks up the curvature of residual moveout for model updates.

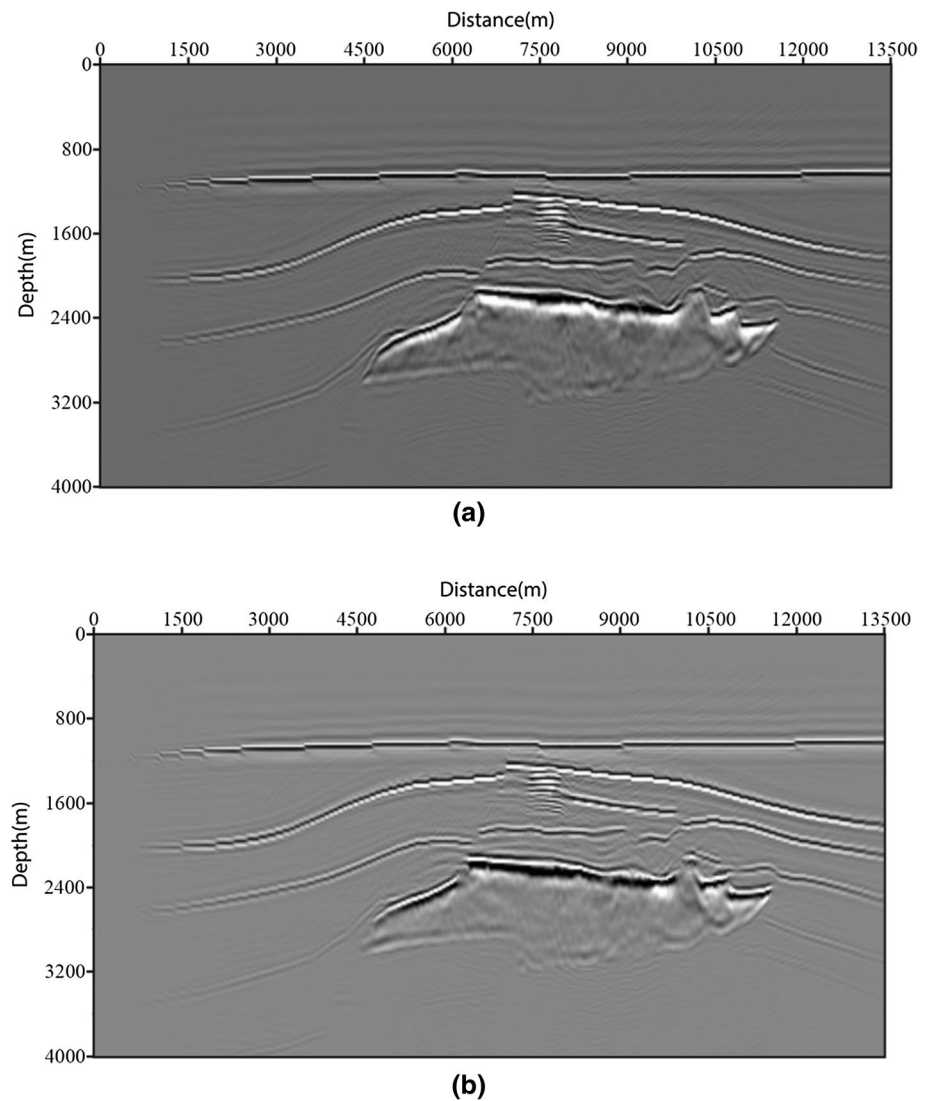
In addition, we compare the computational time and memory costs of the ISO and VTI migration methods shown in Fig. 1. It is obvious that in the conventional 3D VTI migration, the major memory requirement is very large. Even for the current computer cluster, the conventional 3D VTI migration is almost impossible to perform. However, the memory requirement of the proposed PSPI VTI migration method is much lower than the conventional 3D VTI migration and it is almost twice as much as the ISO migration. Compared with the computational efficiency shown in Fig. 1, the proposed PSPI VTI migration method is more efficient than the conventional 3D VTI migration.

The second example illustrates the application of the proposed 3D VTI DSR migration to a marine wide-azimuth seismic dataset from the South China Sea. The field data contain 1300 lines with a 12.5-m line interval and 901 CDPs in each line with a 12.5-m CDP interval. The time extends to 8.192 s with a 4-ms sampling interval. The offset range of hydrophone is  $-6272$  to  $6256$  m in the inline (CDP) direction and  $-1634$  to  $3070$  m in the crossline direction. The frequency range used in the migration process is 2–60 Hz.

The inline section at line 4500 of the velocity model and Thomsen's anisotropic parameters ( $\epsilon$  and  $\delta$ ) is shown in Fig. 7. The model contains mainly anisotropic (VTI) sedimentary layers. In the sedimentary layers,  $\epsilon$  values range from 0 to 0.27 and  $\delta$  values range from 0 to 0.13. There is no significant lateral variation in the horizontal layers of the anisotropic values. Figure 8a, b is the corresponding inline migrated profile of the ISO and proposed VTI migrations, respectively. Comparing the ISO and VTI results, the reflector is corrected by the VTI migration and the seismic event is more continuous and focuses better than the ISO migration. The imaging resolution and signal-to-noise ratio are dramatically enhanced after the VTI migration. This illustrates the validity of the proposed method.

Figure 9 shows three different locations azimuth-opening ADCIGs by the proposed VTI migration, which contain four azimuths with  $90^\circ$  increments and the reflection angle is from  $0^\circ$  to  $60^\circ$  with  $2^\circ$  increment. As we can see, the distribution of angle gathers in different azimuths is different, which reflects the illumination information with azimuthal

**Fig. 5** Inline profile of the migration results generated using **a** ISO DSR extrapolation operator, and **b** proposed VTI DSR extrapolation operator



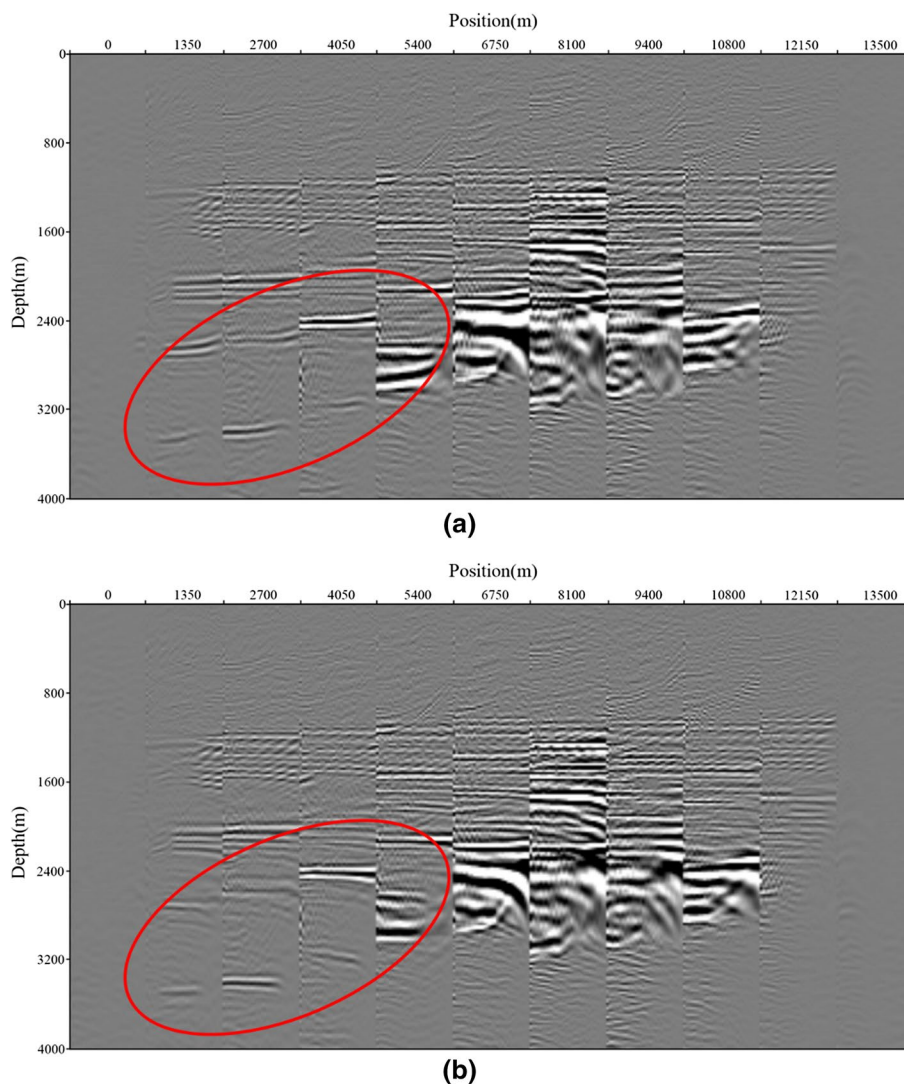
change. Therefore, according to the characteristics of the image gathers, the corresponding illumination information in different directions can be extracted. In addition, we can see that the ADCIGs are well formed and have high resolution. The subsurface angle-dependent reflectivity can be extracted easily. Moreover, these azimuth-opening ADCIGs can be superimposed by optimizing the angle and azimuth to improve the final image quality. Figure 10 shows the angle gathers stacked over all azimuth angle.

The process of preprocessing the migrated data and the image gathers output require very little computation time. Because the MPI and parallel I/O are adopted, the preprocessing of almost 2 Tb data can be completed in about 2 h. Meanwhile, the image gathers calculations can be completed in half an hour. The DSR migration process is time-consuming, and the computational time is related to computing nodes. In our example, we use 30 nodes and each node has 2 processes. (The computer configuration consists of Intel(R)

Xeon(R) CPU E5-2680 v2 @ 2.80 GHz, 40 threads, GeForce GTX 780 Ti in one nodes. In the depth extrapolation, the multi-threaded FFTW with 4 threads and the OPENMP with 8 threads are used in one process.) There are 477 frequency slices, which take 56 h altogether. If we use the same calculation nodes for the frequency slices, it only takes about 8 h to realize the pre-stack depth migration processing.

In the migration process, the RAM and the hard disk storage are concerning issues. The DSR preprocessing part does not consume RAM at all. Because of the new storage approach, the space of the hard disk is largely reduced. In this processing, the size of the prepared data is only 249G. In the migration imaging process, the total RAM is about 3–4 times of the single-frequency data in the current calculated range. Because the image gathers are necessary outputs, the size of the required hard disk depends mainly on the offset range of output. In the post-processing of migration, the RAM is almost not consumed. The output size of the hard disk is related to

**Fig. 6** Inline profile of the common-image gathers generated using **a** ISO DSR extrapolation operator, **b** proposed VTI DSR extrapolation operator



the number of azimuths. It can be determined according to the actual situation. In our example, only four azimuths are output.

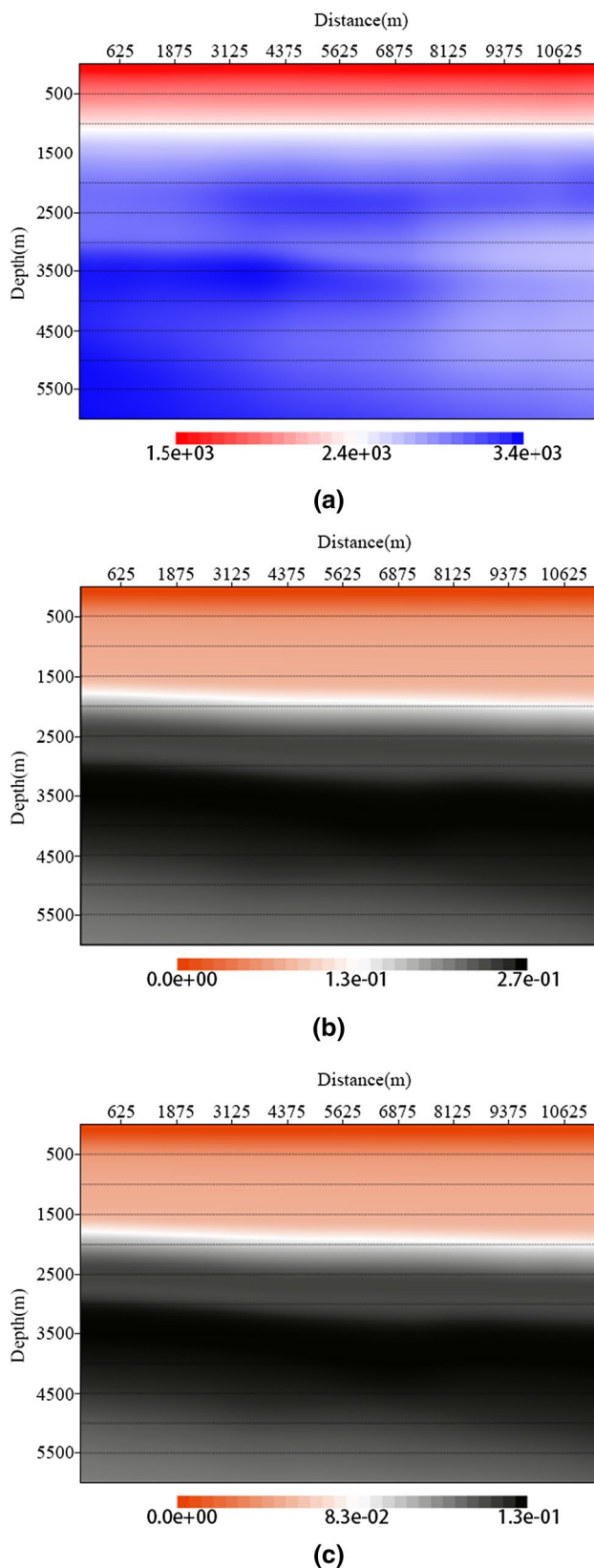
## Discussion

Conventional 3D DSR migration based on narrow azimuth is no longer enough for the current wide-azimuth seismic data. It results in non-negligible errors in the crossline direction. Therefore, the full 3D DSR migration must be considered. In order to accurately describe the wave propagation in anisotropic media, the DSR extrapolation operator must be implemented using Eq. (8). However, the computational cost of the full 3D DSR migration is very expensive and not applicable to large-scale seismic data, so it is necessary to approximate the anisotropic operator. The proposed phase-shift plus interpolation method is a simple approximation method. Its main purpose is to improve the computational efficiency and reduce the memory requirements. If we want

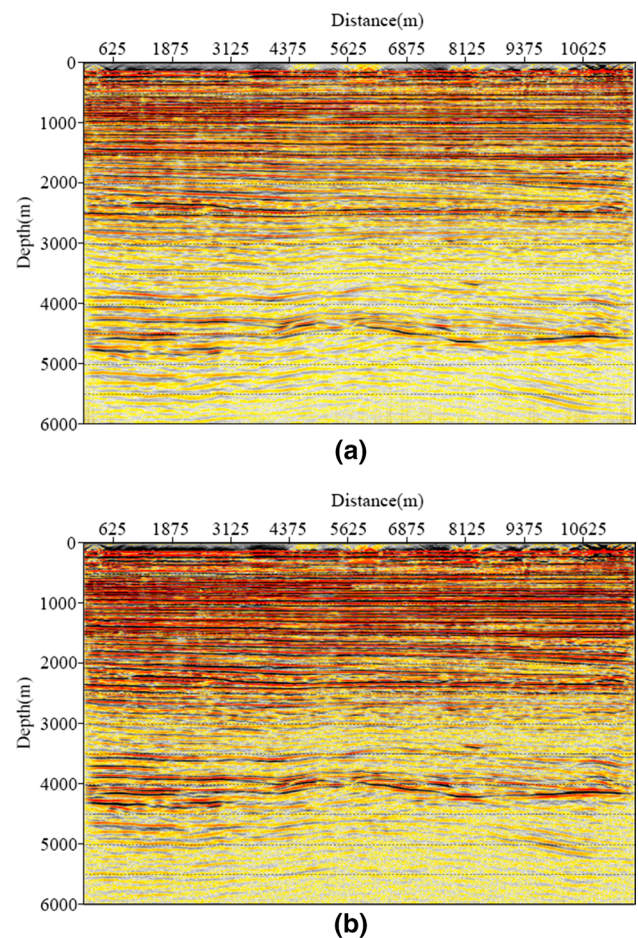
to further improve the image quality, a more precise approximation method needs to be considered. For example, the split domain method (Stoffa et al. 1990) is even more efficient and possibly accurate.

With the different approximations in anisotropic media, we can construct various implementations. For the weak-anisotropy approximation (Thomsen 1986), the interpolation method needs to consider two parameters  $\epsilon$  and  $\delta$  at the same time. For the purpose of improving the computational efficiency, the interpolated coefficient just depends on the parameter  $\epsilon$ . For Alkhalifah's parameterization (Alkhalifah 1998, 2000a), only one parameter  $\eta$  needs to be considered in the interpolation process. In this case, the interpolation is more accurate than the weak-anisotropy approximation.

The proposed angle-domain DSR migration scheme is suitable for any type of seismic data and can be easily matched with the current observation system. The advantage of this strategy is that there is no need to sort the seismic data and generate new data, which could occupy a large



**Fig. 7** Inline profile of marine wide-azimuth seismic dataset from the South China Sea, for **a** migrated velocity model and **b** Thomsen's anisotropic parameters ( $\epsilon$ ), and **c** Thomsen's anisotropic parameters ( $\delta$ )



**Fig. 8** Inline profile of the migration results generated using **a** ISO DSR extrapolation operator, and **b** proposed VTI DSR extrapolation operator

number of hard disks. For a specific computer workstation, only the storage usage of the local and global disks needs to be adjusted. Therefore, the proposed scheme can be conveniently used in real seismic data. In addition, the current 3D angle-domain DSR migration framework is suitable for both isotropic and heterogeneous media. Only the extrapolation operator needs to be modified. The different interpolation methods can be integrated in this framework.

For the azimuth-opening ADCIGs, the number of output azimuth and reflection angles is related to the input observation system and the actual application. The smaller the number, the smaller the memory and storage needed, but the less information available; if the number is large, the pressure of the computer's memory will be prominent. When generating the subsurface migration gathers in the DSR migration step, and converting the subsurface migration gathers into azimuth-opening ADCIGs in the image gathers output step, proper regularization and other processing techniques can significantly improve the quality of image gathers.

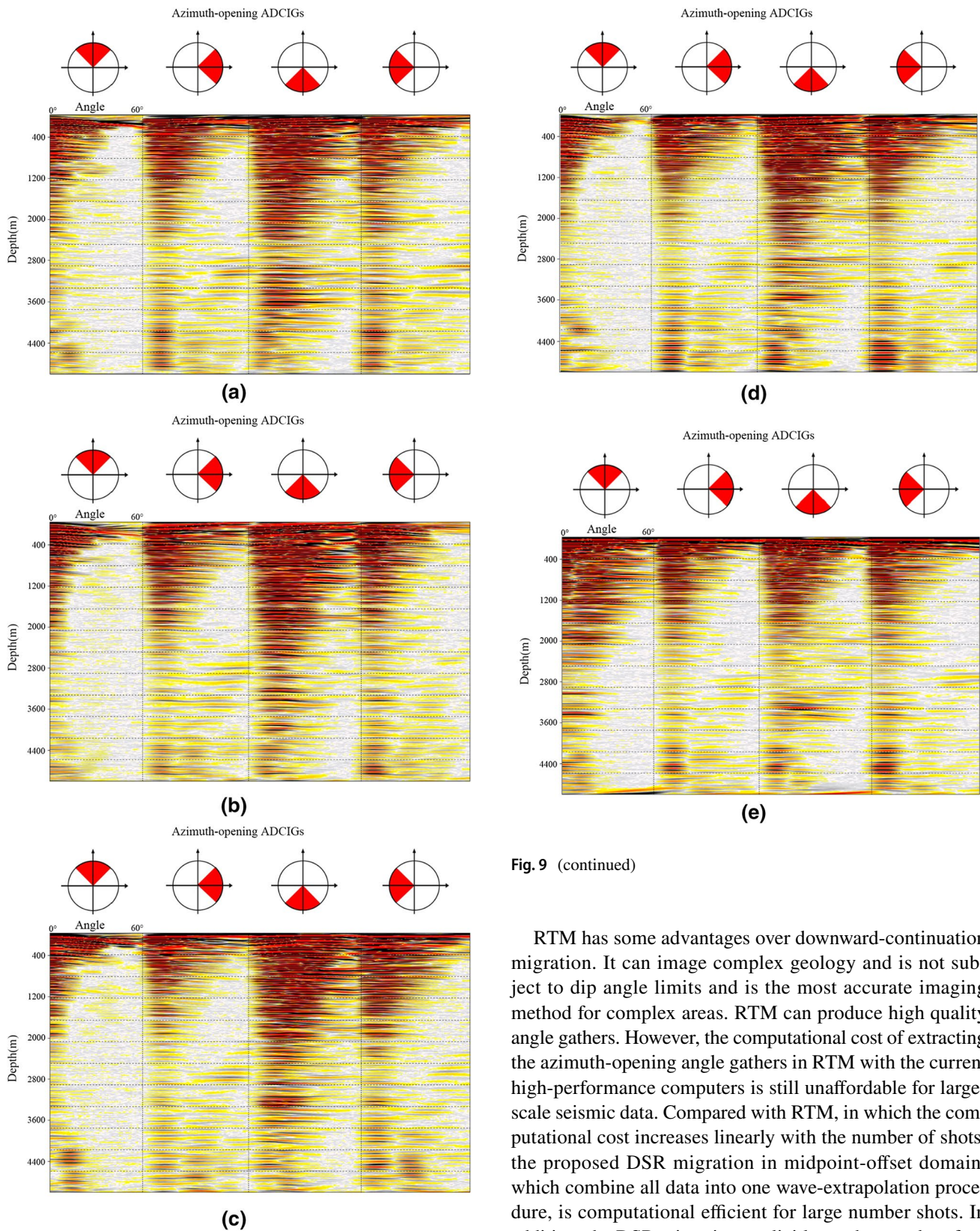
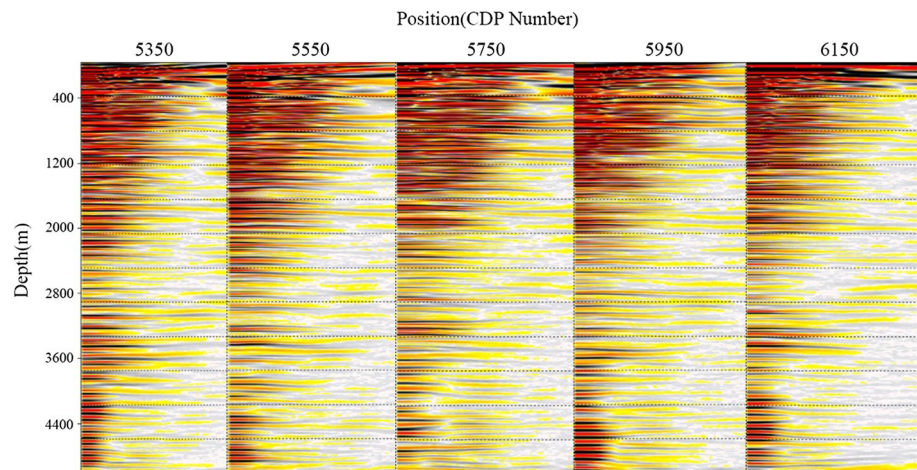


Fig. 9 (continued)

RTM has some advantages over downward-continuation migration. It can image complex geology and is not subject to dip angle limits and is the most accurate imaging method for complex areas. RTM can produce high quality angle gathers. However, the computational cost of extracting the azimuth-opening angle gathers in RTM with the current high-performance computers is still unaffordable for large-scale seismic data. Compared with RTM, in which the computational cost increases linearly with the number of shots, the proposed DSR migration in midpoint-offset domain, which combine all data into one wave-extrapolation procedure, is computational efficient for large number shots. In addition, the DSR migration explicitly produces subsurface offset as part of the wavefield extrapolation, which is convenient in extracting azimuth-opening angle gathers directly. The memory requirement for RTM-ADCIGs is related to the

**Fig. 9** Azimuth-opening ADCIGs generated by the proposed VTI DSR migration. The ADCIGs is selected at five CDP positions (**a** 5350; **b** 5550; **c** 5750; **d** 5950; **e** 6150) and in every ADCIGs, there are four azimuths with  $90^\circ$  increment, and the opening angle range is  $0^\circ$ – $60^\circ$  with a  $1^\circ$  interval

**Fig. 10** ADCIGs by stacking all four azimuths in the corresponding CDP positions



source–receiver aperture, while the memory requirements of DSR are proportional to the offset. Thus, the memory requirement for DSR is larger than that for RTM. Therefore, large memory capacity in the computing node is required for DSR migration.

Moreover, the proposed angle-domain DSR migration can be conveniently applied on the tomographic inversion. In the tomographic inversion, the key factor is whether the azimuth-opening angle gathers can be generated quickly and accurately in the migration process and can be conveniently integrated in migration velocity analysis (MVA). Because migration velocity analysis requires multiple iterations, the migration algorithm also needs to be calculated many times. Therefore, the computational efficiency of the migration algorithm determines whether the MVA can be used in practical applications. With the aid of the high computational efficiency of the proposed method, it is possible for the corresponding MVA to be fast and efficient in practical applications.

The DSR migration is implemented in the frequency domain. The proposed angle-domain DSR migration method can be easily extended to viscoelastic medium, just by changing the frequency to a complex one. As a result, we can apply various attenuation models to realize the amplitude attenuation and phase dispersion and correct for them.

## Conclusions

We propose a full 3D DSR operator for VTI media, which can simultaneously downward extrapolate the whole dataset, efficiently. To realize the wavefield extrapolation, we use the phase-shift plus interpolation method to improve the computational efficiency. Moreover, we show that the 3D azimuth-opening ADCIGs can be generated conveniently in a post-migration scheme, making it possible to

update the velocity and anisotropic parameters through tomographic inversion.

Numerical examples on the salt dome model and real data show that the proposed angle-domain 3D DSR migration can produce high-quality imaging results and output azimuth-opening angle gathers, efficiently. Compared with RTM, the proposed 3D DSR migration method needs large computer memory. The memory requirement is still affordable for modern computer clusters. In addition, the fast ADCIGs outputting scheme are very attractive for 3D model building, which will be published in a companion paper.

**Acknowledgements** The authors thank the sponsors of WPI group for their financial support and help. WPI's research works are also financially supported by National Key R&D Program of China (2019YFC0312004, 2018YFA0702503), National Natural Science Foundation of China (41774126), the great and special project (2016ZX05024-001, 2016ZX05006-002). We kindly acknowledge FFTW Free Software for the computation. We really appreciate the associate editor and Tariq Alkhalifah and an anonymous reviewer for providing so many useful comments and suggestions to improve the clarity and completeness of this manuscript.

## References

- Alkhalifah T (1998) Acoustic approximations for processing in transversely isotropic media. *Geophysics* 63(2):623–631
- Alkhalifah T (2000a) An acoustic wave equation for anisotropic media. *Geophysics* 65(4):1239–1250. <https://doi.org/10.1190/1.1444815>
- Alkhalifah T (2000b) Prestack phase-shift migration of separate offsets. *Geophysics* 65(4):1179–1194
- Alkhalifah T (2000c) The offset-midpoint traveltime equation for transversely isotropic media. *Geophysics* 65(4):1316–1325
- Alkhalifah T (2004) Azimuth moveout correction for transversely isotropic media. *Geophys Prospect* 52(1):39–48
- Alkhalifah T (2012) Prestack exploding reflector modeling and migration for anisotropic media: a parameter estimation tool.

- In 82th SEG annual international meeting. Expanded abstracts, pp 1–5
- Alkhalifah T (2013) Prestack wavefield approximations. *Geophysics* 78(5):T141–T149
- Alkhalifah T (2015) Prestack exploding reflector modeling and migration for anisotropic media. *Geophys Prospect* 63(1):2–10
- Alkhalifah T, Biondi B (2004) Numerical analysis of the azimuth moveout operator for vertically inhomogeneous media. *Geophysics* 69(2):554–561
- Alkhalifah T, Fomel S (2011) Angle gathers in wave-equation imaging for transversely isotropic media. *Geophys Prospect* 59(3):422–431
- Alkhalifah T, Larner K (1994) Migration errors in transversely isotropic media. *Geophysics* 59(9):1405–1418
- Alkhalifah T, Fomel S, Biondi B (2001) The space–time domain: theory and modelling for anisotropic media. *Geophys J Int* 144(1):105–113
- Alkhalifah T, Fomel S, Wu Z (2015) Source-receiver two-way wave extrapolation for prestack exploding-reflector modeling and migration. *Geophys Prospect* 63:23–34
- Barley B, Summers T (2007) Multi-azimuth and wide-azimuth seismic: shallow to deep water, exploration to production. *Lead Edge* 26(4):450–458. <https://doi.org/10.1190/1.2723209>
- Baysal E, Kosloff DD, Sherwood JWC (1983) Reverse time migration. *Geophysics* 48(11):1514–1524. <https://doi.org/10.1190/1.1441434>
- Bevc D, Fliedner M, Crawley S, Biondi B (2003) Wave equation imaging comparisons: survey sinking vs. shot profile methods. In 73th SEG annual international meeting. Expanded abstracts, pp 885–888. <https://doi.org/10.1190/1.1818082>
- Biondi B (2002) Reverse time migration in midpoint-offset coordinates. *SEP Rep* 111:149–156
- Biondi B (2007) Angle-domain common-image gathers from anisotropic migration. *Geophysics* 72(2):581–591
- Biondi B, Chemingui N (1994) Transformation of 3-D prestack data by azimuth moveout (AMO). In 64th SEG annual international meeting. Expanded abstracts, pp 1541–1544. <https://doi.org/10.1190/1.1822833>
- Biondi B, Palacharla G (1996) 3D prestack migration of common-azimuth data. *Geophysics* 61(6):1822–1832. <https://doi.org/10.1190/1.1822730>
- Biondi B, Symes W (2004) Angle-domain common-image gathers for migration velocity analysis by wavefield-continuation imaging. *Geophysics* 69(5):1283–1298. <https://doi.org/10.1190/1.1801945>
- Biondi B, Fomel S, Chemingui N (1998) Azimuth moveout for 3-D prestack imaging. *Geophysics* 63(2):574–588. <https://doi.org/10.1190/1.1444357>
- Bouska J (2008) Advantages of wide-patch, wide-azimuth ocean-bottom seismic reservoir surveillance. *Lead Edge* 27(12):1662–1681. <https://doi.org/10.1190/1.3036972>
- Cai J, Fang W, Wang H (2013) Azimuth–opening angle domain imaging in 3D Gaussian beam depth migration. *J Geophys Eng* 10(2):025013
- Cheng J, Wang H, Ma Z, Yang S (2003) Cross-line common-offset migration for narrow azimuth dataset. In 73th SEG annual international meeting. Expanded abstracts, pp 901–904. <https://doi.org/10.1190/1.1818087>
- Cheng J, Wang H, Ma Z (2005) Double square root equation migration methods of narrow azimuth seismic data. *Chin J Geophys* 48(2):399–405 (in Chinese)
- Cheng J, Ma Z, Geng J, Wang H (2008) Double-square-root one-way wave equation prestack tau migration in heterogeneous media. *Geophys Prospect* 56(1):69–85. <https://doi.org/10.1111/j.1365-2478.2007.00667.x>
- Claerhout JF (1985) *Imaging the Earth's interior*. Blackwell, Oxford
- Cohen JK (1997) Analytic study of the effective parameters for determination of the NMO velocity function in transversely isotropic media. *Geophysics* 62(6):1855–1866
- de Hoop MV, Le Rousseau JH, Wu R-S (2000) Generalization of the phase-screen approximation for the scattering of acoustic waves. *Wave Motion* 31(1):43–70. [https://doi.org/10.1016/S0165-2125\(99\)00026-8](https://doi.org/10.1016/S0165-2125(99)00026-8)
- de Hoop MV, Le Rousseau JH, Biondi B (2003) Symplectic structure of wave-equation imaging: a path-integral approach based on the double-square-root equation. *Geophys J Int* 153(1):52–74. <https://doi.org/10.1046/j.1365-246X.2003.01877.x>
- Duchkov A, de Hoop MV (2009) Extended isochron rays in prestack depth migration. In 79th SEG annual international meeting. Expanded abstracts, pp 3610–3614
- Fomel S (2004) Theory of 3-D angle gathers in wave-equation imaging. In 74th SEG annual international meeting. Expanded abstracts, pp 1053–1056. <https://doi.org/10.1190/1.1851067>
- Frigo M, Steven GJ (2005) The design and implementation of FFTW3. *Proc IEEE* 93(2):216–231
- Gazdag J, Sguazzero P (1984) Migration of seismic data by phase shift plus interpolation. *Geophysics* 49(2):124–131. <https://doi.org/10.1190/1.1441643>
- Gray SH, May WP (1994) Kirchhoff migration using eikonal equation traveltimes. *Geophysics* 59(5):810–817. <https://doi.org/10.1190/1.1443639>
- Hao Q, Stovas A, Alkhalifah T (2015) The offset-midpoint traveltime pyramid in 3D HTI media. *Geophysics* 80(1):T51–T62
- Hao Q, Stovas A, Alkhalifah T (2016) The offset-midpoint traveltime pyramid of P-waves in homogeneous orthorhombic media. *Geophysics* 81(5):C151–C162
- Herman J, Larner K (1995) Prestack migration error in transversely isotropic media. In 65th SEG annual international meeting. Expanded abstracts, pp 1204–1207
- Hu J, Wang H, Wang X (2015) Angle gathers from reverse time migration using analytic wavefield propagation and decomposition in the time domain. *Geophysics* 81(1):S1–S9
- Jin S, Wu R-S (1999) Common offset pseudo-screen depth migration. In 69th SEG annual international meeting. Expanded abstracts, pp 1516–1519. <https://doi.org/10.1190/1.1820809>
- Jin S, Mosher CC, Wu R-S (2002) Offset-domain pseudoscreen prestack depth migration. *Geophysics* 67(6):1895–1902. <https://doi.org/10.1190/1.1527089>
- Jin H, McMechan GA, Guan H (2014) Comparison of methods for extracting ADCIGs from RTM. *Geophysics* 79(3):S89–S103
- Ke B, Zhao B, Liu C, Fang Y (2004) Wave-equation datuming based on a single shot gather. In 74th SEG annual international meeting. Expanded abstracts, pp 2156–2159. <https://doi.org/10.1190/1.1845210>
- Li V, Guitton A, Tsvankin I, Alkhalifah T (2018) Image-domain wavefield tomography for VTI media. *Geophysics* 84(2):1MA–Z11
- Liu Z, Bleistein N (1995) Migration velocity analysis: theory and an iterative algorithm. *Geophysics* 60(1):142–153. <https://doi.org/10.1190/1.1443741>
- Liu S, Wang H, Yang Q (2014) Traveltime computation and imaging from rugged topography in 3D TTI media. *J Geophys Eng* 11(1):0150031
- Liu S, Wang H, Feng B (2015) The characteristic wave decomposition and imaging in VTI media. *J Appl Geophys* 115:51–58
- Liu S, Gu H, Tang Y, Han B, Wang H, Liu D (2018) Angle-domain common imaging gather extraction via Kirchhoff prestack depth migration based on a traveltime table in transversely isotropic media. *J Geophys Eng* 15(2):568–575
- Michell S, Shoshitaishvili E, Chergotis D, Sharp J, Etgen J (2006) Wide azimuth streamer imaging of Mad Dog; Have we solved the subsalt imaging problem? In 76th SEG annual international



- meeting. Expanded abstracts, pp 2905–2909. <https://doi.org/10.1190/1.2370130>
- Mulder WA, Plessix RE (2003) One-way and two-way wave-equation migration. In 73rd SEG annual international meeting. Expanded abstracts, pp 881–884. <https://doi.org/10.1190/1.1818081>
- Oh JW, Alkhalifah T (2018) Optimal full-waveform inversion strategy for marine data in azimuthally rotated elastic orthorhombic media. *Geophysics* 83(4):R307–R320
- Popovici AM (1996) Prestack migration by split-step DSR. *Geophysics* 61(5):1412–1416
- Reshet M (1991) Depth migration from irregular surfaces with depth extrapolation methods. *Geophysics* 56(1):119–122. <https://doi.org/10.1190/1.1442947>
- Sava P, Alkhalifah T (2013) Wide-azimuth angle gathers for anisotropic wave-equation migration. *Geophys Prospect* 61(1):75–91. <https://doi.org/10.1111/j.1365-2478.2012.01024.x>
- Sava PC, Fomel S (2003) Angle-domain common-image gathers by wavefield continuation methods. *Geophysics* 68(3):1065–1074. <https://doi.org/10.1190/1.1581078>
- Sava P, Fomel S (2005) Coordinate-independent angle-gathers or wave equation migration. In 75th SEG annual international meeting. Expanded abstracts, pp 2052–2055
- Sava P, Vlad I (2011) Wide azimuth angle gathers for wave equation migration. *Geophysics* 76(3):S131–S141
- Sava PC, Biondi B, Fomel S (2001) Amplitude-preserved common image gathers by wave-equation migration. In 71st SEG annual international meeting. Expanded abstracts, pp 296–299. <https://doi.org/10.1190/1.1816598>
- Song X, Fomel S (2011) Fourier finite-difference wave propagation. *Geophysics* 76(5):T123–T129. <https://doi.org/10.1190/geo2010-0287.1>
- Stoffa PL, Fokkema JT, de Luna Freire RM, Kessinger WP (1990) Split-step Fourier migration. *Geophysics* 55(4):410–421
- Sun H, Huang L, Fehler MC (2005) Globally optimized Fourier finite-difference migration in the offset domain. In 75th SEG annual international meeting. Expanded abstracts, pp 1858–1861. <https://doi.org/10.1190/1.2148065>
- Thomsen L (1986) Weak elastic anisotropy. *Geophysics* 51(10):1954–1966. <https://doi.org/10.1190/1.1442051>
- Tsvankin I (2012) Seismic signatures and analysis of reflection data in anisotropic media. Society of Exploration Geophysicists, 3rd ed
- VerWest BJ, Lin D (2007) Modeling the impact of wide-azimuth acquisition on subsalt imaging. *Geophysics* 72(5):SM241–SM250. <https://doi.org/10.1190/1.2736516>
- Whitmore ND (1983) Iterative depth migration by backward time propagation. In 53rd SEG annual international meeting. Expanded abstracts, pp 382–385. <https://doi.org/10.1190/1.1893867>
- Wu R-S (1994) Wide-angle elastic wave one-way propagation in heterogeneous media and an elastic wave complex-screen method. *J Geophys Res* 99:751–766. <https://doi.org/10.1029/93JB02518>
- Wu R-S (1996) Synthetic seismograms in heterogeneous media by one-return approximation. *Pure Appl Geophys* 148:155–173
- Wu G, Liang K, Wang H (2007) High-order one-way generalized-screen propagation operator of qP-wave in VTI medium. *Oil Geophys Prospect* 42(6):640–650 (in Chinese)
- Wu C, Wang H, Zhou Y, Hu J (2019) Angle-domain common-image gathers in reverse-time migration by combining the Poynting vector with local-wavefield decomposition. *Geophys Prospect* 67(8):2035–2060
- Xu S, Zhang Y, Tang B (2011) 3D angle gathers from reverse time migration. *Geophysics* 76(2):S77–S92
- Yan R, Xie XB (2012) AVA analysis based on RTM angle domain common image gather. In 82nd SEG annual international meeting. Expanded abstracts, pp 1–6. <https://doi.org/10.1190/segam2012-0521.1>
- Yan L, Larry RL, Lawton DC (2004) Influence of seismic anisotropy on prestack depth migration. *Geophysics* 23(1):30–36
- Yilmaz O, Claerbout JF (1980) Prestack partial migration. *Geophysics* 45(12):1753–1779
- Yuan SY, Wang SX, Luo YN, Wei WW, Wang GC (2019) Impedance inversion by using the low-frequency full-waveform inversion result as an a priori model. *Geophysics* 84(2):R149–R164



# Stable absorption compensation with lateral constraint

Xiong Ma<sup>1</sup> · Guofa Li<sup>1</sup> · Hao Li<sup>1</sup> · Jiaojiao Li<sup>1</sup> · Xu Fan<sup>2</sup>

Received: 12 December 2019 / Accepted: 9 June 2020 / Published online: 26 June 2020  
© Institute of Geophysics, Polish Academy of Sciences & Polish Academy of Sciences 2020

## Abstract

The presence of seismic absorption distorts seismic record and reduces seismogram resolution, which can be partially compensated by application of absorption compensation algorithms. Conventional absorption compensation techniques are based on 1D forward model with each seismic trace being compensated independently. Therefore, the 2D results combined by each compensation trace may be noisy and discontinuity. To eliminate this issue, we extend the 1D forward model to the 2D forward system and further add an additional lateral constraint to the compensation algorithm for enforcing the lateral continuity of the compensated section. Solving the proposed laterally constrained absorption compensation (LCAC) problem, we simultaneously obtain the multiple compensated traces with lateral smoother transition and higher signal-to-noise ratio ( $S/N$ ). We testify the effectiveness of the proposed method by applying both synthetic and field data. Synthetic data examples demonstrate the superior performance of the LCAC algorithm in terms of improving algorithmic stability and protecting lateral continuity. The field data tests further indicate its ability to not only improve seismic resolution, but also inhibit the amplification of high-frequency noise.

**Keywords** Absorption compensation · Lateral constraint · Seismic resolution · Lateral continuity

## Introduction

Seismic wave propagating through the Earth undergoes seismic absorption due to the anelasticity of the subsurface medium, resulting in the reduction in vertical resolution and the stretching of seismic wavelets (Kolsky 1956; Futterman 1962; Kjartansson 1979). The quality factor  $Q$  is commonly used to quantitatively characterize these absorption effects (Li et al. 2016), and it is also a basic parameter in the absorption compensation processing (Wang 2002; Zhang and Ulrych 2007). So far, various absorption compensation schemes, including nonstationary deconvolution (Margrave et al. 2011; Van der Baan 2012; Oliveira

and Lupinacci 2013; Yuan et al. 2017), inverse  $Q$  filtering (Robinson 1979; Bickel and Natarajan 1985; Hargreaves and Calvert 1991; Wang 2006; Li et al. 2015; Wang et al. 2018b), and  $Q$ -compensated migration (Mittet et al. 1995; Dutta and Schuster 2014; Zhao et al. 2018; Wang et al. 2018c, 2019), have been developed to compensate for the absorption of seismic energy and to enhance the resolution of seismic data.

Inverse  $Q$  filtering, also known as absorption compensation, has been addressed by many researchers (Zhang and Ulrych 2007; Braga and Moraes 2013; Chai et al. 2014). The core problem of inverse  $Q$  filtering is the inherent instability of the amplitude compensation, which includes an exponential amplification term in the compensation operator (Zhang and Ulrych 2007). In recent years, many attempts have been explored to suppress the high-frequency noise amplification and further to obtain a stable compensation solution (Wang 2002; Zhang and Ulrych 2007; Braga and Moraes 2013; Chai et al. 2014; Yuan et al. 2016). To our knowledge, the stabilized strategies can be classified into two categories. The first category focuses on modifying the compensation operator to achieve a stable result. For example, Robinson (1979) proposes a phase-only inverse  $Q$  filtering, which neglects the amplitude exponential amplification term, to correct the phase distortion in the seismic

✉ Guofa Li  
lgfseismic@126.com

Xiong Ma  
mx\_geophysics@126.com

<sup>1</sup> China University of Petroleum-Beijing, School of Geophysics, State Key Lab of Petroleum Resources and Prospecting, Key Lab of Geophysical Exploration of CNPC, Changping, Beijing 102249, China

<sup>2</sup> Institute of Geophysics, Institute of Exploration and Development, Xinjiang Oilfield, Ürümqi 830000, Xinjiang, China

data. Wang (2002) presents a gain-limited inverse  $Q$  filtering, which implements amplitude compensation only within the limited frequency band, to avoid amplifying of high-frequency noise. He further develops the stabilized inverse  $Q$  filtering by adding a regularization factor into the amplitude compensation operator to improve the stability of absorption compensation. The second category formulates the absorption compensation as an inverse problem by minimizing the misfit between observed data and modeled data (Chai et al. 2014; Wang et al. 2018b). In the objective function, they only exploit the forward (exponential decay) operator and do not require the inverse exponential amplification operator, which means the instability due to exponential amplitude amplification is avoided by using inverse scheme. However, as discussed by Wang (2011), the compensation result by using inverse scheme is about solving the first kind Fredholm integral equation; thus, its numerical solution is unstable. For obtaining a stable solution, we need to incorporate some prior information or constraints in the inversion framework. For example, Zhang and Ulrych (2007) use the Cauchy–Gauss prior model to regularize the inverse problem by means of Bayes' theorem. Braga and Moraes (2013) apply the  $L_2$  norm constraint to accomplish inverse  $Q$  filtering in the wavelet domain. Wang et al. (2018b) present a  $L_{1-2}$ -regularized absorption compensation algorithm for stable seismic compensation. Nevertheless, the above absorption compensation methods are based on 1D forward model and apply trace-by-trace compensation strategy; thus, the compensated 2D section combined by each 1D result may be noisy and shows a poor lateral continuity (Auken and Christiansen 2004; Auker et al. 2005; Hamid and Pidlisecky 2015; Wang et al. 2018a; Ma et al. 2019; Ji et al. 2019; Yuan et al. 2019; Ma et al. 2020).

To reduce the lateral discontinuity problems in the trace-by-trace inversion algorithm, Auker and Christiansen (2004) originally develop a laterally constrained inversion algorithm for resistivity data processing. Afterward, Schmalz and Tezkan (2007) use this algorithm for transit electromagnetic inversion. Hamid and Pidlisecky (2015) further introduce it to seismic exploration and develop a lateral constraint algorithm for seismic impedance inversion. In addition, Zhang et al. (2013) modify the lateral constraint to the 'Z' shape constraint for multi-trace seismic reflectivity inversion. In this paper, we incorporate the lateral constraint between adjacent seismic traces into the absorption compensation processing and furthermore present a laterally constrained absorption compensation (LCAC) algorithm to enforce the lateral continuity of the compensated section. Synthetic and field data examples indicate that the proposed LCAC method improves seismic resolution and lateral continuity.

The structure of this paper is as follows: Firstly, we briefly review the conventional 1D absorption compensation algorithm. Then, we extend the 1D algorithm to the

2D algorithm and incorporate a lateral constraint term into inversion system for developing a novel LCAC algorithm. Next, we use synthetic and field data experiments to verify the superiority of the proposed LCAC method in terms of improving algorithmic stability and protecting lateral continuity. Finally, we draw some conclusions.

## Theory and method

### Laterally unconstrained absorption compensation

In the elastic medium, the post-stack seismic record can be modeled by the convolution of a source wavelet with the reflectivity sequences (Yilmaz 2001),

$$s_0(t) = w(t) \otimes r(\tau), \quad (1)$$

where the notation  $\otimes$  represents the convolutional operator,  $t$  and  $\tau$  are the record time,  $s_0(t)$  is the non-attenuated seismic trace,  $w(t)$  is the source wavelet, and  $r(\tau)$  is the reflectivity series.

When considering seismic wave propagation in absorption medium, Eq. 1 can be modified as the nonstationary convolution model (Margrave 1998):

$$s(t) = \hat{w}(t, \tau) \otimes r(\tau), \quad (2)$$

where  $s(t)$  is the attenuated seismic trace and  $\hat{w}(t, \tau)$  is the time-varying wavelet due to  $Q$  filtering effects which can be expressed as,

$$\hat{w}(t, \tau) = \int_0^\infty W(\omega) A(\omega, \tau) e^{i\omega t} d\omega, \quad (3)$$

where  $\omega$  is angular frequency,  $i$  denotes the imaginary unit,  $W(\omega)$  is the frequency spectrum of the source wavelet  $w(t)$ , and  $A(\omega, \tau)$  is the  $Q$ -filtering function determined by the selected absorption model. In this paper, we use the modified Kolsky–Futterman model (Wang and Guo 2004) to describe seismic wave propagation in absorption media. Then, the  $Q$ -filtering function  $A(\omega, \tau)$  is expressed as:

$$A(\omega, \tau) = \exp \left[ -i\omega\tau \left| \frac{\omega_r}{\omega} \right|^{-\gamma} \left( 1 - \frac{i}{2Q} \right) \right], \quad (4)$$

where  $\omega_r$  represents the reference angular frequency,  $\gamma = \frac{1}{\pi Q}$  is a dimensionless factor, and  $Q$  is the quality factor.

According to Eqs. (1)–(4), we can derive the relationship between the attenuated signal  $s(t)$  and the non-attenuated signal  $s_0(t)$  which is written by (see Appendix and Braga and Moraes 2013):

$$s(t) = a(t, \tau) \otimes s_0(t), \quad (5)$$

where  $a(t, \tau) = \int_0^\infty A(\omega, \tau)e^{i\omega t}d\omega$  is the time domain  $Q$ -filtering function.

Equation 5 is the basis of 1D absorption compensation, and it can be interpreted that the attenuated seismogram is obtained by the convolution of the attenuation function with the non-attenuated signal. The matrix-vector form of Eq. 5 is,

$$\mathbf{s} = \mathbf{A}\mathbf{s}_0, \tag{6}$$

where  $\mathbf{s}$  and  $\mathbf{s}_0$  are, respectively, the vectors of the attenuated and non-attenuated seismic signals and the matrix

$$\mathbf{A} = \begin{bmatrix} a(t_1, \tau_1) & a(t_1, \tau_2) & \cdots & a(t_1, \tau_N) \\ a(t_2, \tau_1) & a(t_2, \tau_2) & \cdots & a(t_2, \tau_N) \\ \vdots & \vdots & \ddots & \vdots \\ a(t_N, \tau_1) & a(t_N, \tau_2) & \cdots & a(t_N, \tau_N) \end{bmatrix}$$

stands for the  $Q$  filtering effects.

According to Braga and Moraes (2013), the numerical solution  $\mathbf{s}_0$  of Eq. 6 is unstable and we should apply some regularization terms to stabilize the compensation (non-attenuated) result. By using the  $L_2$  norm regularization, the following objective functional is established which can be represented by,

$$\mathcal{J}(\mathbf{s}_0) = \|\mathbf{A}\mathbf{s}_0 - \mathbf{s}\|^2 + \lambda\|\mathbf{s}_0\|^2, \tag{7}$$

where  $\|\bullet\|^2$  represents  $L_2$  norm and  $\lambda$  is the regularization parameter.

Equation 7 is a standard least-squares problem, and its solution can be expressed as (Braga and Moraes 2013):

$$\mathbf{s}_0 = \mathbf{M}^{-1}\mathbf{b}, \tag{8}$$

where  $\mathbf{M} = \mathbf{A}^T\mathbf{A} + \lambda\mathbf{I}$ ,  $\mathbf{b} = \mathbf{A}^T\mathbf{s}$  and the superscript  $T$  denotes the transpose. Using Eq. 8, Braga and Moraes (2013) compensate a 2D seismic section trace-by-trace and then combine all 1D results to form a 2D compensation section. In this algorithm, he only regularize the inverted solution in the vertical (or time) direction but with the lateral direction unconstrained, so we refer it to as laterally unconstrained absorption compensation (LUAC) algorithm.

In the LUAC algorithm, the regularization parameter  $\lambda$  has some influences on the compensation result. The parameter  $\lambda$  should be small when the  $S/N$  of seismic data is high and  $\lambda$  should be relatively larger when the  $S/N$  of seismic data is lower. For determining a suitable  $\lambda$ , we can use L-curve technique (Hansen and O’Leary 1993), which applies a cross-plot of the data error versus the solution length as a function of  $\lambda$ . A good value for the parameter is the one located at the corner of the L-curve.

### Laterally constrained absorption compensation

The LUAC algorithm neglects the lateral constraint in its objective functional (Eq. 7); thus, the compensated profile

may be discontinuous in the lateral direction when the attenuated signals are contaminated by random noise. Therefore, in this section, we focus on incorporating a lateral constraint term into the inversion objective functional and further developing a LCAC algorithm.

For taking the lateral continuity information into consideration, we should extend the 1D forward model (Eq. 6) to multichannel forward system (Ma et al. 2020):

$$\mathbf{d} = \mathbf{G}\mathbf{d}_0, \tag{9}$$

where  $\mathbf{d} = [\mathbf{s}_1, \mathbf{s}_2, \dots, \mathbf{s}_M]^T$  and  $\mathbf{d}_0 = [\mathbf{s}_{01}, \mathbf{s}_{02}, \dots, \mathbf{s}_{0M}]^T$  are, respectively, the concatenated attenuated data vector and

$$\mathbf{G} = \begin{bmatrix} \mathbf{A}_1 & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{A}_2 & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \ddots & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{A}_M \end{bmatrix}$$

is a non-attenuated data vector, and  $\mathbf{G}$  is a

block diagonal matrix representing the multichannel  $Q$  filtering effects.

Based on multichannel forward model (Eq. 9), we take the lateral constraint into consideration and set up a laterally constrained objective functional:

$$\mathcal{J}(\mathbf{d}_0) = \|\mathbf{G}\mathbf{d}_0 - \mathbf{d}\|^2 + \lambda\|\mathbf{d}_0\|^2 + \mu\|\mathbf{D}_x\mathbf{d}_0\|^2, \tag{10}$$

where  $\mu$  is the lateral regularization parameter, which controls the relative strength of the lateral constraint term to the data misfit term, and

$$\mathbf{D}_x = \begin{bmatrix} -1 & 0 & \cdots & 0 & 1 & 0 & \cdots & 0 & 0 & 0 \\ 0 & -1 & 0 & \cdots & 0 & 1 & 0 & \cdots & 0 & 0 \\ \vdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \vdots & \cdots \\ 0 & 0 & 0 & \cdots & 0 & -1 & 0 & \cdots & 0 & 1 \end{bmatrix}$$

is the horizontal first-order derivative matrix.

The least-squares solution of this problem is:

$$\mathbf{d}_0 = \tilde{\mathbf{M}}^{-1}\tilde{\mathbf{b}}, \tag{11}$$

where  $\tilde{\mathbf{M}} = \mathbf{G}^T\mathbf{G} + \lambda\mathbf{I} + \mu\mathbf{D}_x^T\mathbf{D}_x$  and  $\tilde{\mathbf{b}} = \mathbf{G}^T\mathbf{d}$ . Compared with the LUAC method, the proposed LCAC algorithm compensates all seismic traces simultaneously and protects the lateral continuity of the inverted results. The overall performance of the LCAC approach is verified by using synthetic and field data examples in the next section.

For the proposed LCAC method, there are two regularization parameters,  $\lambda$  and  $\mu$ , to select and these two parameters regularize the strength of vertical and lateral constraints, respectively. Until now, the optimization for hyperparameters functional remains a complex problem (Clapp et al. 2004). In this paper, we select them by trial and error, but we apply a relatively elegant strategy. Firstly, we set the lateral regularization parameter  $\mu = 0$  and use L-curve technique to choose a proper parameter  $\lambda$ . Secondly, we fix the parameter  $\lambda$  calculated in the first step and then optimize  $\mu$  by trial and error. For synthetic data test, we can evaluate the compensated results by comparing them with referenced data based

on similarity or cross-correlation criterion. For field data application, we can judge the compensation results from two aspects, that is, seismic resolution and  $S/N$ . A good result has not only improved seismic resolution but also higher  $S/N$ .

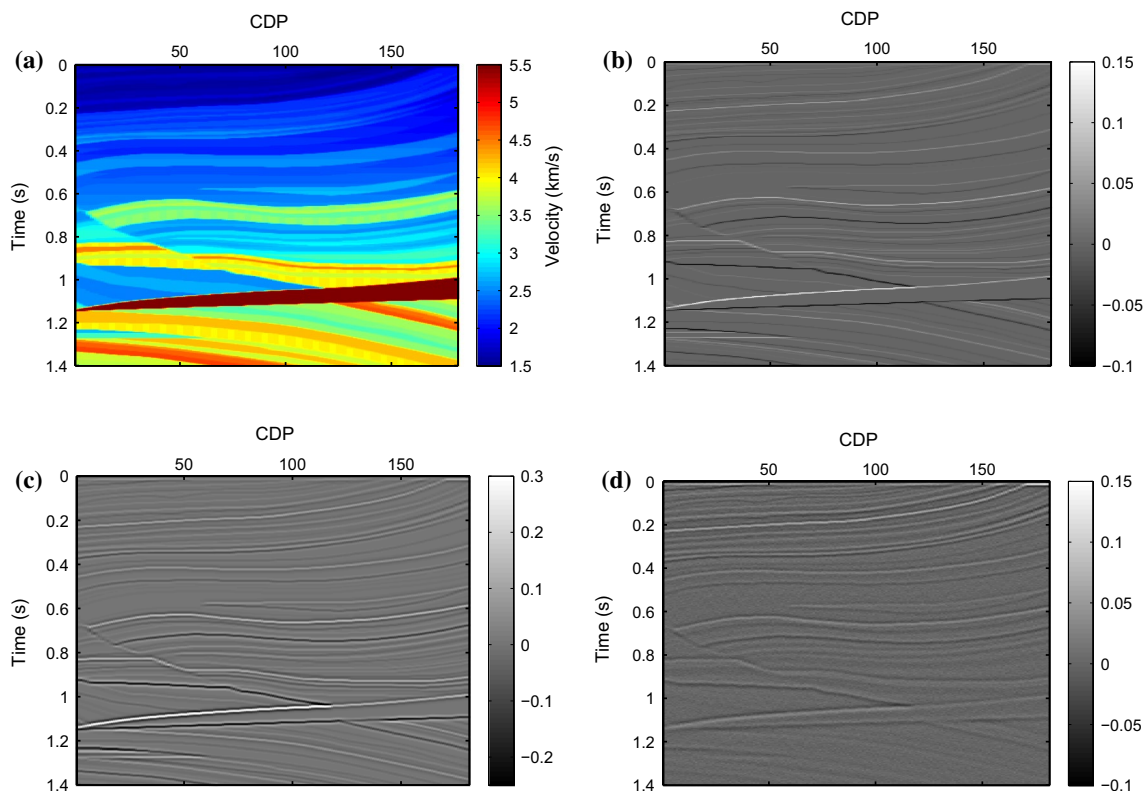
## Examples

### Synthetic data example

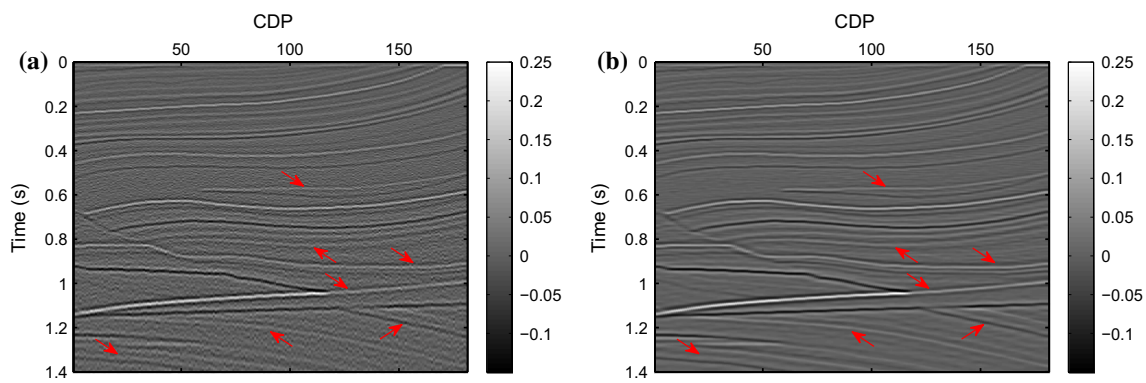
In the section, we exploit the partial Marmousi model to demonstrate the effectiveness and superiority of the proposed LCAC algorithm. Figure 1a shows the velocity model. We assume the density is constant, and then, we obtain the reflectivity model displayed in Fig. 1b. Figure 1c displays the stationary seismic data generated by convolving a 40 Hz Ricker wavelet with the above reflectivity model. The non-attenuated data are served as the reference data to evaluate the compensation performance of both LUAC and LCAC approaches. Figure 2d depicts the attenuated seismogram with the quality factor  $Q = 50$  and contaminated by 20% Gaussian noise. Due to the  $Q$  filtering effects, the energy of deep reflection events is attenuated and the resolution of recorded seismic data is decreased.

We use both LUAC and LCAC algorithms to process the attenuated seismic data for recovering the seismic events and improving the seismic resolution. In both methods, we apply the true  $Q$  value as input. In LUAC method, we choose the parameter  $\lambda = 0.008$  and display the inverted section in Fig. 2a. As we can see, the LUAC compensated section partially recovers the seismic reflections and enhances the vertical resolution of seismic data, but the lateral continuity of compensated data is poor and the  $S/N$  is low. Moreover, we calculate the correlation coefficient (CC) between the LUAC result and the reference data (Fig. 1c) and the value is 0.7566. In the proposed LCAC algorithm, we fix  $\lambda = 0.008$  and determine the lateral regularization parameter as  $\mu = 0.5$ . The corresponding compensation data are shown in Fig. 2b. We observe that the LCAC algorithm generates a result with better lateral continuity (see red arrows) and with relatively higher  $S/N$  than LUAC result. The CC between it and reference data reaches to 0.8902.

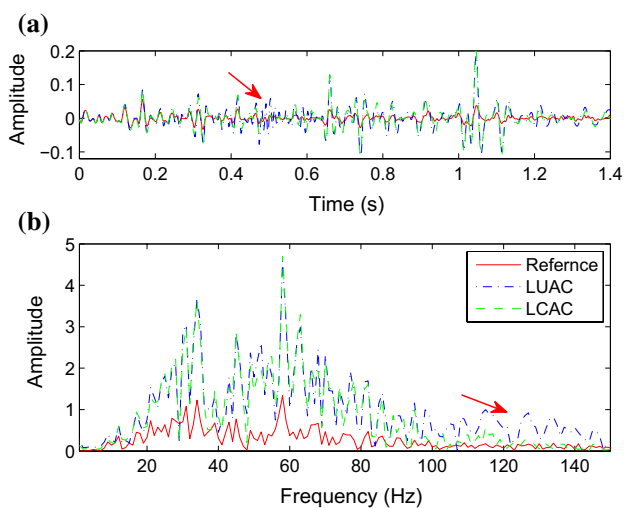
For a clear comparison, seismic traces extracted from the attenuated data (Fig. 1d) and the compensated results (Fig. 2) at CDP=101 are shown in Fig. 3a and their corresponding spectra are displayed in Fig. 3b. From the extracted traces, we see that the overall compensation performance of both algorithms is similar except for the seismic noise amplification in some places (see red arrow). The comparison of



**Fig. 1** Forward modeling for generating stationary and nonstationary data. **a** The velocity model, **b** the reflectivity model, **c** the synthetic stationary seismic data without seismic noise, and **d** the synthetic nonstationary (attenuated) data with 20% Gaussian noise



**Fig. 2** Absorption compensation results. **a** The LUAC compensation data with the regularization parameter  $\lambda = 0.008$ , **b** the LCAC compensation result with the vertical regularization parameter  $\lambda = 0.008$ , and the lateral regularization parameter  $\mu = 0.5$



**Fig. 3** Attenuated and compensated seismic traces and their spectra. **a** The seismic traces extracted from Figs. 1d and 2 at CDP=101, and **b** their spectra

their spectra further confirms that the proposed method can not only recover seismic events but also suppress the high-frequency noise amplification.

We also use the attenuated data shown in Fig. 1d to study the influence of the lateral regularization parameter  $\mu$  on the compensation results. We fix the parameters  $\lambda = 0.005$ , and, respectively, select  $\mu$  as 5, 0.5, and 0.01. The corresponding compensated data are displayed in Fig. 4a–c, respectively. When  $\mu$  is too large, the compensated data appear to be over-smoothed and the faults are blurry (see arrows). When  $\mu$  is too small (Fig. 4c), the amplification of seismic noise is evident. When  $\mu$  is moderate (Fig. 4b), the compensated results achieve a good balance between the noise suppression and the lateral continuity enhancement. Figure 5 shows the attenuated seismic traces, the reference traces, and the LCAC compensated traces at CDP=100. It is confirmed that

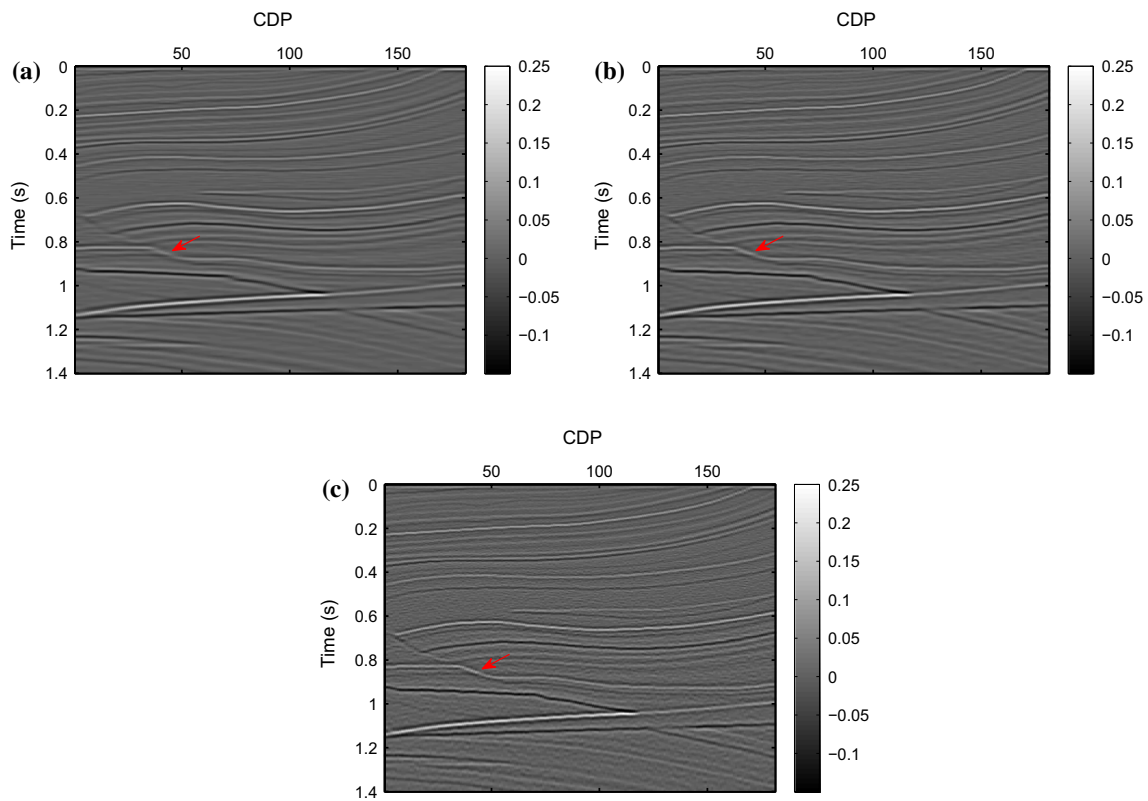
the large lateral regularization parameter  $\mu$  leads to the result (Fig. 5a) with stronger noise suppression (see green arrows) and with relatively weaker signal recovery (see blue arrows). Therefore, we should be careful about enforcing the lateral continuity too strong because sometimes the spatial discontinuity could be a response of the real geological structure, such as faults or pinch-out.

Noticing that the quality factor  $Q$  is difficult to estimate, and it is important to examine the effects of using inaccurate  $Q$  values in the proposed LCAC algorithm. The attenuated data shown in Fig. 1d are exploited again to conduct the experiments. In the tests, the true  $Q$  value is  $Q_{\text{True}} = 50$ . The relative error of  $Q$  factor is defined as:

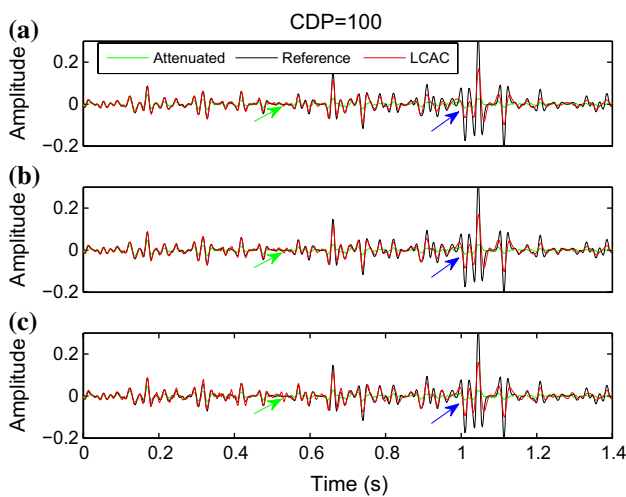
$$e_r = \frac{Q - Q_{\text{true}}}{Q_{\text{true}}} \times 100\%. \quad (12)$$

Moreover, we use the CC of the compensated data and the reference data (Fig. 1c) to evaluate the compensation performance.

Table 1 shows the correlation coefficients by using several different  $Q$  values. Since the input data (Fig. 1d) contain random noise, the LCAC results cannot be very close to the reference data (Fig. 1c) even if the  $Q$  value is accurate. Specifically, when using an accurate  $Q$  value  $Q = 50$ , the CC is 0.8902 (shown in row 3 of Table 1) and it is treated as a reference value for our comparison below. As expected, when using the inaccurate  $Q$  values, the correlation coefficients of the LCAC results are decreased, which means the precision of the compensation results is reduced by using inaccurate  $Q$  values. In the case of a slightly inaccurate  $Q$  (e.g.,  $Q = 45$  or  $55$ ), the correlation coefficients (e.g.,  $CC = 0.8536$  or  $0.8783$ ) are close to the reference value  $CC = 0.8902$ . This means a slightly inaccurate  $Q$  value produces only minor perturbations in the LCAC results. In the case of a moderately inaccurate  $Q$  (e.g.,  $Q = 40$  or  $60$ ), although the corresponding



**Fig. 4** Investigation of the influence of the lateral regularization parameter  $\mu$  on the compensation results. We fix the vertical regularization parameter  $\lambda = 0.005$  and, respectively, select  $\mu$  as **a** 5, **b** 0.5, and **c** 0.01



**Fig. 5** Comparisons of the compensation performance with difference lateral regularization parameters. **a**  $\mu = 5$ , **b**  $\mu = 0.5$ , and **c**  $\mu = 0.01$

correlation coefficients (e.g.,  $CC = 0.8142$  or  $0.8586$ ) are decreased compared with those using slightly inaccurate  $Q$  values, these compensation results are still acceptable. Furthermore, as shown in Table 1, overestimation of  $Q$

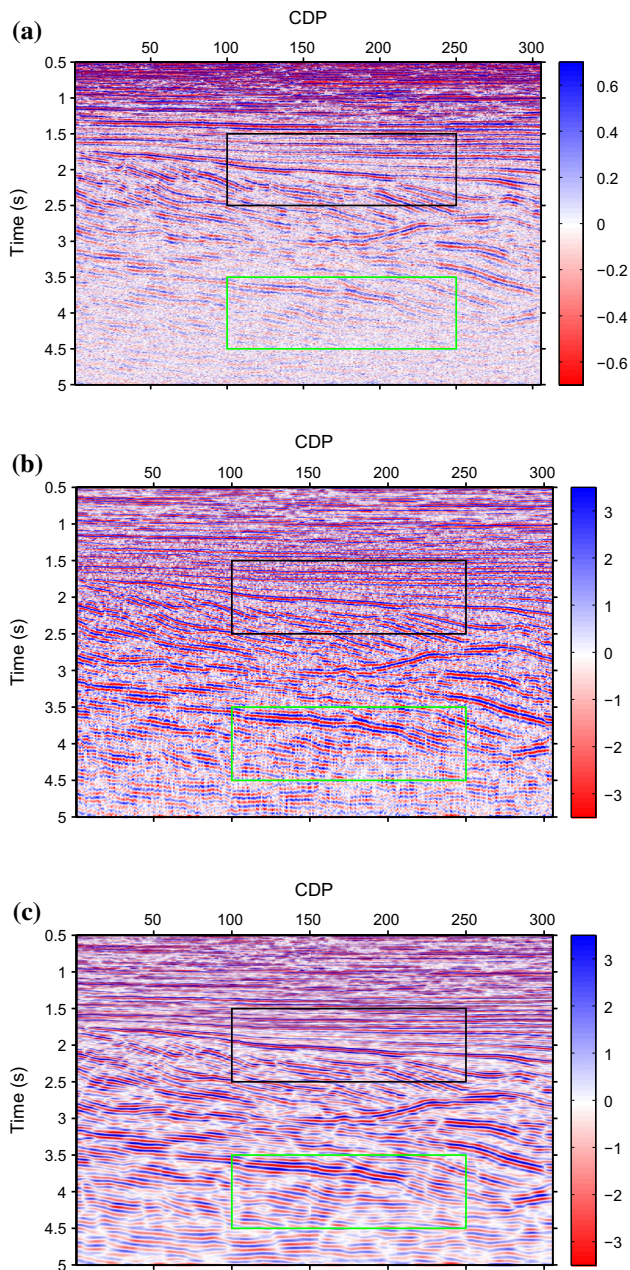
**Table 1** Correlation coefficients by using different  $Q$  values

No.	$Q_{true}$	$Q$	$e_r$	CC
1	50	40	-20%	0.8142
2	50	45	-10%	0.8536
3	50	50	0%	0.8902
4	50	55	10%	0.8783
5	50	60	20%	0.8586
6	50	100	100%	0.8077

(the positive  $Q$  error) may have less influence on the compensation results than that underestimation of  $Q$  (the negative  $Q$  error). A too small  $Q$  value will overcompensate for seismic absorption, while a large  $Q$  value will lead to undercompensation. In the absorption compensation, we prefer undercompensating for seismic absorption to overcompensating so as to suppressing the high-frequency noise amplification. In the tests, the relative error of  $Q$  ranges from  $-20\%$  to  $100\%$ , but the CC of the LCAC results shows relatively small fluctuation, which indicates that relatively imprecise  $Q$  estimates can still produce acceptable compensation results.

## Field data tests

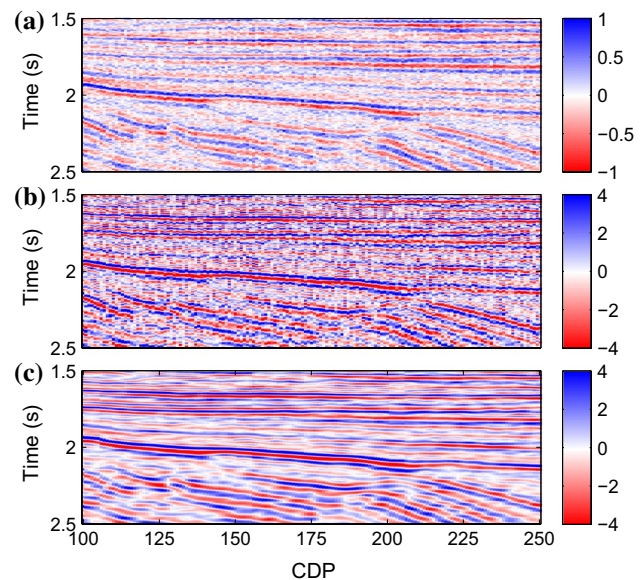
For further verifying the practicability of the LCAC algorithm, we apply both the LUAC and LCAC algorithms to field data as shown in Fig. 6a. This field data are acquired in East China and free of absorption compensation processing. The energy of deep reflections is weak, and the resolution of the raw seismic data is poor. We use both LUAC and LCAC methods to compensate for the raw data.



**Fig. 6** Seismic attenuation compensation for field data. **a** The field data, **b** the compensation result from LUAC method, and **c** the compensation result from LCAC approach

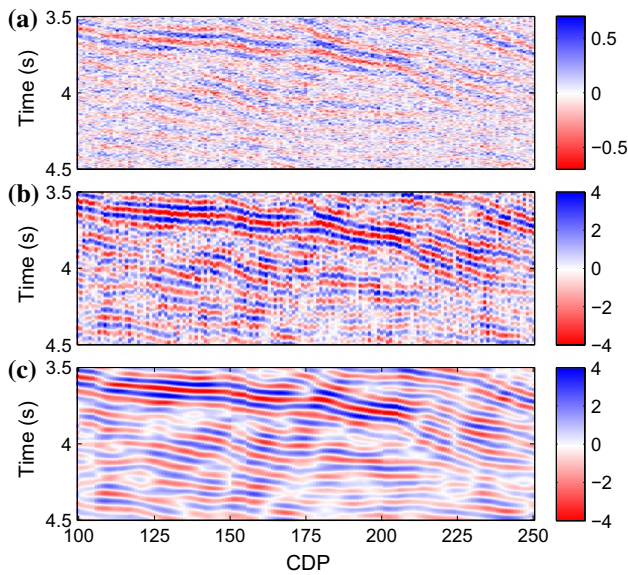
Before implementing absorption compensation, we estimate  $Q$  value via the attenuation-based  $Q$  analysis (Ma et al. 2017). Figure 6b, c display the compensation sections by using LUAC and LCAC methods respectively. Compared with the raw data, two compensated results partially recover the seismic data absorption and enhance the seismic resolution. The further comparison of LUAC and LCAC compensation sections indicates that the LCAC algorithm provides a result with higher  $S/N$  and smoother spatial continuity without losing evident vertical resolution.

For viewing more detailed compensation features, we display the black and green boxes portion of Fig. 6 in a zoomed view (Figs. 7, 8). From these figures, we find that the proposed LUAC compensation results inhibit the high-frequency noise amplification and exhibit an improved lateral continuity compared with the LUAC result, which demonstrates that the proposed method is more robust to random noise. Figure 9 shows the amplitude spectra of the raw data, the LUAC compensation result, and the LCAC compensation result. We observe that the amplitude spectra of compensated data are broadened and the mid-high frequency components are boosted after absorption compensation processing, but the LUAC algorithm has boosted more seismic energy owing to the amplification of seismic noise.

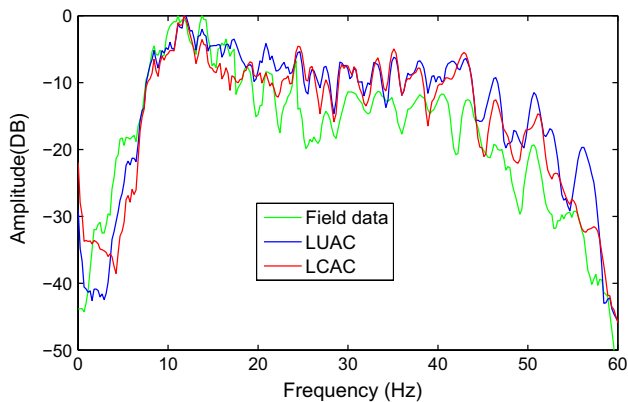


**Fig. 7** Zoomed view of the black boxes shown in Fig. 6. **a** The raw data, **b** the LUAC compensation result, and **c** the LCAC compensation result





**Fig. 8** Zoomed view of the green boxes shown in Fig. 7. **a** The raw data, **b** the LUAC compensation result, and **c** the LCAC compensation result



**Fig. 9** Amplitude spectra of the raw data, the LUAC result, and the LCAC result

## Discussion

In recent years, many multichannel approaches have been developed to deal with the nonstationary seismic data, such as nonstationary multichannel reflectivity inversion or nonstationary multichannel deconvolution (Haghshenas Lari and Gholami 2019). The objective function of these methods can be written in the following formula,

$$\mathcal{J}(\mathbf{x}) = \|\mathbf{K}\mathbf{x} - \mathbf{y}\|^2 + \lambda V(\mathbf{x}) + \mu H(\mathbf{x}), \quad (13)$$

where  $\mathbf{x}$  is the desired result,  $\mathbf{y}$  is the input data,  $\mathbf{M}$  is the operator linking  $\mathbf{x}$  to  $\mathbf{y}$ , and  $V(\mathbf{x})$  and  $H(\mathbf{x})$  are, respectively, the vertical (time) and horizontal constraints.

Compared with the multichannel deconvolution methods, there are certain differences between them and our proposed approach. Firstly, the output result  $\mathbf{x}$  is different in these two kind of approaches, that is, the desired result  $\mathbf{x}$  is reflectivity sequences in the multichannel deconvolution algorithms (Haghshenas Lari and Gholami 2019), while it is the non-attenuated or compensated seismic data in our multichannel compensation algorithm (see Eq. 10). The reason for this difference lies in the different goals. For the stationary multichannel deconvolution algorithms, the main goal is to remove the wavelet effects and extract the high-resolution reflectivity sequences from the stationary seismic data (Du et al. 2018). If the input seismic data is nonstationary, the stationary deconvolution algorithms can be extended to the nonstationary deconvolution algorithms which may simultaneously eliminate the wavelet-filtering and  $Q$ -filtering effects (Haghshenas Lari and Gholami 2019). This means the nonstationary deconvolution algorithms have certain theoretical advantages. However, as we all know, both wavelet estimation and  $Q$ -compensation are great challenges in field data processing. Therefore, the nonstationary deconvolution algorithms which try to deal with these two problems simultaneously may show relatively poor applicability in field data processing. In this paper, we treat the wavelet estimation (or elimination) and  $Q$ -compensation as two separated problems and the proposed multichannel compensation algorithm focuses on compensating for seismic absorption due to  $Q$ -filtering effects. Thus, we do not need any wavelet information in our algorithm. This also explains why the output result  $\mathbf{x}$  of the proposed method is non-attenuated or compensated seismic records rather than reflectivity sequences. After we obtain the compensated seismic data, many stationary deconvolution algorithms can be used to further get the high-resolution reflectivity sequences. This two-step strategy has relatively strong applicability in field data application.

Secondly, the proposed method has higher computational efficiency than the multichannel deconvolution algorithms. In the multichannel deconvolution algorithms, a sparse regularization, e.g.,  $L_1$  norm regularization, is usually imposed on the reflectivity sequences. In other words, the vertical constraint  $V(\mathbf{x})$  in Eq. 13 can be written as  $V(\mathbf{x}) = \|\mathbf{x}\|_1$ . Because the  $L_1$  norm regularization is a nonlinear function, the objective function of the deconvolution algorithms is also nonlinear. To solve this nonlinear problem, many iterative algorithms, such as iterative reweighting algorithm (Sacchi 1997), split the Bregman method (Haghshenas Lari and Gholami 2019) and alternating direction method of multipliers (Du et al. 2018), are generally used which transform the nonlinear problem into linear problem in each iteration

step. The computational cost depends on the number of iterations and the cost in each iteration step. As for the proposed algorithm, we employ the  $L_2$  norm in both the vertical and horizontal constraints (see Eq. 10); thus, we obtain a linear objective function which can be directly solved by employing gradient algorithms, such as conjugate gradient method.

## Conclusions

In this paper, we incorporate the lateral constraint into absorption compensation algorithm and develop a LCAC method. Compared with LUAC approach, the proposed LCAC method provides a compensation profile with better lateral continuity and higher  $S/N$ , which may be more geologically realistic. The synthetic data tests demonstrate the strong stability of the proposed LCAC method in enhancing the vertical resolution while suppressing the seismic noise amplification. The application of field data further indicates its practicability and viability as a robust absorption compensation method. In addition, we should be careful about enforcing the lateral regularization too strong because sometimes the lateral discontinuity may be a response of the real geological structure. In our future research, we will focus on improving our algorithm to be more edge-preserving.

**Acknowledgements** We would like to acknowledge the financial support by National Key R & D Program of China (Grant No. 2018YFA0702504) and National Natural Science Foundation of China (Grant No. 41874141).

## Compliance with ethical standards

**Conflict of interest** On behalf of all authors, the corresponding author states that there is no conflict of interest.

## Appendix: Derivation of Eq. 5

As we all know, the convolution in the time domain is equal to the product in the frequency domain, and then, the frequency domain expressions of Eq. (1) and Eq. (2) are

$$S_0(\omega) = W(\omega)R(\omega), \quad (\text{A.1})$$

and

$$S(\omega) = \hat{W}(\omega, \tau)R(\omega), \quad (\text{A.2})$$

where  $S_0(\omega)$  and  $S(\omega)$  are, respectively, the non-attenuated and attenuated seismic signals in the frequency domain,  $R(\omega)$  is the frequency domain reflectivity, and  $\hat{W}(\omega, \tau)$  is the time-varying wavelet in the frequency domain. According to Eq. (3), the frequency domain time-varying wavelet can be written by:

$$\hat{W}(\omega, \tau) = W(\omega)A(\omega, \tau), \quad (\text{A.3})$$

Substituting Eq. (A.3) back into Eq. (A.2), we have,

$$S(\omega) = W(\omega)A(\omega, \tau)R(\omega) = A(\omega, \tau)S_0(\omega), \quad (\text{A.4})$$

By transforming Eq. (A.4) into the time domain, we obtain Eq. (5),

$$s(t) = a(t, \tau) \otimes s_0(t), \quad (\text{A.5})$$

where  $a(t, \tau) = \int_0^\infty A(\omega, \tau)e^{i\omega t} d\omega$  is the inverse Fourier transform of frequency domain  $Q$ -filtering function.

## References

- Auken E, Christiansen AV (2004) Layered and laterally constrained 2D inversion of resistivity data. *Geophysics* 69(6):752–761
- Auken E, Christiansen AV, Jacobsen BH, Foged N, Sørensen KI (2005) Piecewise 1D laterally constrained inversion of resistivity data. *Geophys Prospect* 53(4):497–506
- Bickel SH, Natarajan RR (1985) Plane-wave Q deconvolution. *Geophysics* 50(9):1426–1439
- Braga ILS, Moraes FS (2013) High-resolution gathers by inverse filtering in the wavelet domain. *Geophysics* 78(2):V53–V61
- Chai X, Wang S, Yuan S, Zhao J, Sun L, Wei X (2014) Sparse reflectivity inversion for nonstationary seismic data. *Geophysics* 79(3):V93–V105
- Clapp R, Biondi B, Claerbout JF (2004) Incorporating geologic information into reflection tomography. *Geophysics* 69(2):533–546
- Du X, Li G, Zhang M, Li H, Yang W, Wang W (2018) Multichannel band-controlled deconvolution based on a data-driven structural regularization. *Geophysics* 83(5):R401–R411
- Dutta G, Schuster GT (2014) Attenuation compensation for least-squares reverse time migration using the viscoacoustic-wave equation. *Geophysics* 79(6):S251–S262
- Futterman WI (1962) Dispersive body waves. *J Geophys Res* 67(13):5279–5291
- Haghshenas Lari H, Gholami A (2019) Nonstationary blind deconvolution of seismic records. *Geophysics* 84(1):V1–V9
- Hamid H, Pidlisecky A (2015) Multitrace impedance inversion with lateral constraints. *Geophysics* 80(6):M101–M111
- Hansen P, O’Leary DP (1993) The use of L-curve in the regularization of discrete ill-posed problems. *SIAM J Sci Comput* 14:1487–1503
- Hargreaves ND, Calvert AJ (1991) Inverse Q filtering by Fourier transform. *Geophysics* 56(4):519–527
- Ji Y, Yuan S, Wang S (2019) Multi-trace stochastic sparse-spike inversion for reflectivity. *J Appl Geophys* 161:84–91
- Kjartansson E (1979) Constant Q wave propagation and attenuation. *J Geophys Res Solid Earth* 84(B9):4737–4748
- Kolsky H (1956) LXXI. The propagation of stress pulses in viscoelastic solids. *Philos Mag* 1(8):693–710
- Li G, Liu Y, Zheng H, Huang W (2015) Absorption decomposition and compensation via a two-step scheme. *Geophysics* 80(6):V145–V155
- Li G, Sacchi MD, Zheng H (2016) In situ evidence for frequency dependence of near-surface Q. *Geophys J Int* 204(2):1308–1315
- Ma X, Li G, Wang S, Yang W, Wang W (2017) A new method for Q estimation from reflection seismic data. In: 87th Annual International Meeting, SEG Expanded Abstracts, 5496–5500
- Ma M, Zhang R, Yuan SY (2019) Multichannel impedance inversion for nonstationary seismic data based on the modified alternating direction method of multipliers. *Geophysics* 84(1):A1–A6

- Ma X, Li G, Li H, Yang W (2020) Multichannel absorption compensation with a data-driven structural regularization. *Geophysics* 85(1):V71–V80
- Margrave GF (1998) Theory of nonstationary linear filtering in the Fourier domain with application to time-variant filtering. *Geophysics* 63(1):244–259
- Margrave GF, Lamoureux MP, Henley DC (2011) Gabor deconvolution: estimating reflectivity by nonstationary deconvolution of seismic data. *Geophysics* 76(3):W15–W30
- Mittet R, Sollie R, Hokstad K (1995) Prestack depth migration with compensation for absorption. *J Appl Geophys* 34(2):1485–1494
- Oliveira SAM, Lupinacci WM (2013) L1 norm inversion method for deconvolution in attenuating media. *Geophys Prospect* 61(4):771–777
- Robinson JC (1979) A technique for the continuous representation of dispersion in seismic data. *Geophysics* 44(8):1345–1351
- Sacchi MD (1997) Reweighting strategies in seismic deconvolution. *Geophys J Int* 129:651–656
- Schmalz T, Tezkan B (2007) 1D-Laterally Constraint Inversion (1D-LCI) of Radiomagnetotelluric Data from a Test Site in Denmark. *Kolloquium Elektromagnetische Tiefenforschung* 199–204
- Van der Baan M (2012) Bandwidth enhancement: inverse Q filtering or time-varying Wiener deconvolution? *Geophysics* 77(4):V133–V142
- Wang Y (2002) A stable and efficient approach of inverse Q filtering. *Geophysics* 67(2):657–663
- Wang Y (2006) Inverse Q-filter for seismic resolution enhancement. *Geophysics* 71(3):V51–V60
- Wang SD (2011) Attenuation compensation method based on inversion. *Appl Geophys* 8(2):150–157
- Wang Y, Guo J (2004) Modified Kolsky model for seismic attenuation and dispersion. *J Geophys Eng* 1(3):187–196
- Wang Y, Liu W, Cheng S, She B, Hu G, Liu W (2018a) Sharp and laterally constrained multitrace impedance inversion based on blocky coordinate descent. *Acta Geophys* 66:623–631
- Wang Y, Ma X, Zhou H, Chen Y (2018b)  $L_{1-2}$  minimization for exact and stable seismic attenuation compensation. *Geophys J Int* 213(3):1629–1646
- Wang Y, Zhou H, Chen H, Chen Y (2018c) Adaptive stabilization for Q-compensated reverse time migration. *Geophysics* 83(1):S15–S32
- Wang Y, Zhou H, Zhao X, Zhang Q, Chen Y (2019) Q-compensated viscoelastic reverse time migration using mode-dependent adaptive stabilization scheme. *Geophysics* 84(4):S301–S315
- Yilmaz O (2001) *Seismic data analysis*. Society of Exploration Geophysicists, Tulsa
- Yuan S, Wang S, Tian N, Wang Z (2016) Stable inversion-based multitrace deabsorption method for spatial continuity preservation and weak signal compensation. *Geophysics* 81(3):V199–V212
- Yuan S, Wang S, Ma M, Ji Y, Deng L (2017) Sparse Bayesian learning-based time-variant deconvolution. *IEEE Trans Geosci Remote Sens* 55(11):6182–6194
- Yuan S, Wang S, Luo Y, Wei W, Wang G (2019) Impedance inversion by using the low-frequency full-waveform inversion result as an a priori model. *Geophysics* 84(2):R149–R164
- Zhang C, Ulrych TJ (2007) Seismic absorption compensation: a least squares inverse scheme. *Geophysics* 72(6):R109–R114
- Zhang R, Sen MK, Srinivasan S (2013) Multi-trace basis pursuit inversion with spatial regularization. *J Geophys Eng* 10(3):035012
- Zhao X, Zhou H, Wang Y, Chen H, Zhou Z, Sun P, Zhang J (2018) A stable approach for Q-compensated viscoelastic reverse time migration using excitation amplitude imaging condition. *Geophysics* 83(5):S459–S476



# Three-dimensional magnetotelluric inversion using L-BFGS

Libin Lu<sup>1</sup> · Kunpeng Wang<sup>2</sup> · Handong Tan<sup>1</sup> · Qingkun Li<sup>2</sup>

Received: 11 November 2019 / Accepted: 18 June 2020 / Published online: 30 June 2020  
© Institute of Geophysics, Polish Academy of Sciences & Polish Academy of Sciences 2020

## Abstract

The gradient-based optimization methods are preferable for the large-scale three-dimensional (3D) magnetotelluric (MT) inverse problem. Compared with the popular nonlinear conjugate gradient (NLCG) method, however, the limited-memory Broyden–Fletcher–Goldfarb–Shanno (L-BFGS) method is less adopted. This paper aims to implement a L-BFGS-based inversion algorithm for the 3D MT problem. And we develop our code on top of the ModEM package, which is highly extensible and popular among the MT community. To accelerate the convergence speed, the preconditioning technique by the affine linear transformation of the original model parameters is used. Two modifications of the conventional L-BFGS algorithm are also made to get a comparable convergence rate with the NLCG method. The impacts of the preconditioner parameters, the regularization parameters, the starting model, etc., on the inversion are evaluated by synthetic examples for both L-BFGS and NLCG methods. And the real MT Kayabe dataset is also inverted by the inversion algorithms. The synthetic tests show that through our L-BFGS inversion algorithm the similar resistivity models can be obtained with that from the NLCG method. For the real data inversion, the L-BFGS method performs more efficiently and reasonable results could be obtained by less iterations of the inversion process than the NLCG method. Thus, we suggest the common usage of the L-BFGS method for the 3D MT inverse problem.

**Keywords** 3D · MT · Line search · NLCG · Quasi-Newton

## Introduction

Three-dimensional magnetotelluric inversion has been receiving substantial attention in the context of complex geological structures for the last two decades (Newman et al. 2008; Sass and Ritter et al. 2014; Devi et al. 2019). To solve the inverse problem, various optimization methods could be taken into consideration which are classified into two categories. The sensitivity-based methods, such as the Gauss–Newton method (Jahandari and Farquharson 2017) and data-space Occam method (Siripunvaraporn and Egbert 2000), requiring computing the second-order derivatives of the objective functional, can be time-consuming and resource-intensive despite good convergence properties. When dealing with large-scale model and large

dataset, the gradient-based methods, such as the steepest descent method and nonlinear conjugate gradient method, are of greater interest, especially for NLCG method which has gained great popularity among the MT community due to its clarity and effectiveness (Newman and Alumbaugh 2000; Rodi and Mackie 2001; Kelbert et al. 2008; Kelbert et al. 2014). Siripunvaraporn and Sarakorn (2011) combined the conjugate gradient method with the data-space Occam method, so as to obtain the convergence behavior from the sensitivity information and reduce the computation time and memory. Likewise, the quasi-Newton (QN) method respects the sensitivity information by approximating the Hessian matrix, which could result in similar benefits. The limited-memory Broyden–Fletcher–Goldfarb–Shanno algorithm, belonging to the family of the quasi-Newton method, proved to be a competitive method (Liu and Nocedal 1989). It only calculates and stores a small number of gradient vector difference pairs and model vector difference pairs to approximate the Hessian matrix (Byrd et al. 1994), thus reducing memory requirement. Furthermore, the L-BFGS method can be adjusted to solve the bound constraint problem in which the parameters are limited within a given range (Byrd et al.

✉ Kunpeng Wang  
xfnwkp@163.com

<sup>1</sup> School of Geophysics and Information Technology, China University of Geosciences (Beijing), Beijing 100083, China

<sup>2</sup> College of Geophysics, Chengdu University of Technology, Chengdu 610059, China

1995). And it is very meaningful for a geophysical inverse problem like magnetotellurics.

Provided that the L-BFGS method has shown promising features, it is increasingly employed in geophysical electromagnetic inverse problems. For example, Newman and Boggs (2004) adopted the L-BFGS method in three-dimensional cross-well electromagnetic imaging. As for magnetotellurics, Avdeeva and Avdeev (2006) applied the limited memory quasi-Newton method to one-dimensional magnetotelluric inversion. By applying the QN scheme proposed by Ni and Yuan (1997), their algorithm had the capability of limiting the resistivity values within a given range. They further extended the algorithm to the 3D MT case and tested it on synthetic models (Avdeev and Avdeeva 2009; Avdeeva et al. 2012). Moorkamp et al. (2011) also chose the L-BFGS method to jointly invert multiple kinds of geophysical data, including MT data.

It was shown that using a preconditioner on the NLCC or L-BFGS method could speed up the convergence (Newman and Boggs 2004). The approximate Hessian as such a preconditioner is a possible choice but brings additional computation cost. Avdeev and Avdeeva (2009) introduced an additional regularization technique through which the original gradient was scaled by a sequence of coefficients. Not only could the erratic structures causing by the singularity of the gradient be eliminated, but also less iterations were needed in their example. Apart from those, a comparable preconditioning could be achieved by the affine linear transformation of the model parameters (Kelbert et al. 2008; Egbert and Kelbert 2012). And they implemented this transformation technique with the NLCC method in the ModEM code (Kelbert et al. 2014). Liu and Yin (2013) tried the affine linear transformation for the 3D Helicopter electromagnetic inversion. While the ModEM code is readily in use and absent of the L-BFGS method, we have developed the L-BFGS inversion for 3D MT into which the transformation technique is ported. Another aspect of the inversion is the line search scheme. It is doubtful to run large number of line search involving large-scale problems. To avoid excessive computation cost, we propose a relaxation line search scheme in this paper.

The outline of this paper is as follows. In “Methods” section, we give the methodology behind our approach. The 3D MT inverse problem is briefly expressed, together with the affine linear transformation. Two modifications of the original L-BFGS code, including the proposed line search strategy, are explained. We also state the major differences between the NLCC and L-BFGS method in this section. In “Results and discussion” section, through synthetic and real data tests, we will investigate the impact of the controlling parameters, such as smooth factor and smooth number in constructing the model covariance matrix, on the stability and convergence rate of the inversion. We also compare

the results of the L-BFGS inversion with that by the NLCC inversions. Finally, we draw conclusions in “Conclusions” section.

## Methods

### Objective functional and gradient calculation

The magnetotelluric inversion, together with other geophysical electromagnetic inversions, is to seek a best model by minimizing the following objective functional

$$f(\mathbf{m}) = [\mathbf{d} - \mathbf{F}(\mathbf{m})]^{T*} \mathbf{C}_d^{-1} [\mathbf{d} - \mathbf{F}(\mathbf{m})] + \lambda (\mathbf{m} - \mathbf{m}_0)^T \mathbf{C}_m^{-1} (\mathbf{m} - \mathbf{m}_0) \quad (1)$$

where  $\mathbf{m}$  is the current model,  $\mathbf{d}$  is the measured data, and  $\mathbf{F}(\mathbf{m})$  indicates forward mapping over the given model.  $\lambda$  is the regularization parameter or trade-off parameter which controls the balance of the data fidelity term and regularization term.  $\mathbf{C}_d$  and  $\mathbf{C}_m$  are the data covariance matrix and model covariance matrix, respectively. And  $\mathbf{m}_0$  is a prior model or reference model (Siripunvaraporn and Egbert 2000; Siripunvaraporn et al. 2005).

The minimization of the functional in Eq. 1 can be solved by an iterative optimization method, such as the L-BFGS method. Different from the sensitivity-based method, it only requires the gradient of the functional to get a model update. And the gradient of the functional is given by

$$\nabla f(\mathbf{m}) = -2\text{Re}\{\mathbf{J}^T \mathbf{C}_d^{-1} [\mathbf{d} - \mathbf{F}(\mathbf{m})]^*\} + 2\lambda \mathbf{C}_m^{-1} (\mathbf{m} - \mathbf{m}_0) \quad (2)$$

where  $\mathbf{J}$  is the sensitive matrix. Although the sensitivity matrix is included in Eq. 2, the reciprocal method can be used to avoid forming the explicit sensitivity matrix. Instead, the product of  $\mathbf{J}^T$  with a given vector is computed. In such way, the computation time and memory usage are significantly reduced (Newman and Alumbaugh 2000; Rodi and Mackie 2001). Kelbert et al. (2008) introduced the following affine linear transformation

$$\tilde{\mathbf{m}} = \mathbf{C}_m^{-1/2} (\mathbf{m} - \mathbf{m}_0) \quad (3)$$

and changed the original objective functional to the following form,

$$f(\tilde{\mathbf{m}}) = [\mathbf{d} - \mathbf{F}(\tilde{\mathbf{m}})]^{T*} \mathbf{C}_d^{-1} [\mathbf{d} - \mathbf{F}(\tilde{\mathbf{m}})] / N_d + \lambda \tilde{\mathbf{m}}^T \tilde{\mathbf{m}} / N_m \quad (4)$$

Here  $N_d$  and  $N_m$  are the total numbers of data and model parameters, respectively. By dividing with these two entities, the corresponding data term and model regularization term tend to be dimensionless. Then the gradient of the new objective functional is

$$\nabla f(\tilde{\mathbf{m}}) = -2\mathbf{C}_m^{1/2} \text{Re}\{\mathbf{J}^T \mathbf{C}_d^{-1} [\mathbf{d} - \mathbf{F}(\tilde{\mathbf{m}})]^*\} / N_d + 2\lambda \tilde{\mathbf{m}} / N_m \tag{5}$$

And the original model parameters can be recovered as follows:

$$\mathbf{m} = \mathbf{C}_m^{1/2} \tilde{\mathbf{m}} + \mathbf{m}_0 \tag{6}$$

From Eq. 5, we see that  $\mathbf{C}^{1/2} \mathbf{m}$  serves as a smooth operator on the original gradient. Kelbert et al. (2008) suggested it as a preconditioner for the NLCG inversion. So, we assume it can be also incorporated into the L-BFGS method.

### The formation of the smooth operator

Among most geophysical inversion framework, the formation of the model covariance matrix  $\mathbf{C}_m$  or its square root  $\mathbf{C}_m^{1/2}$  would result from a finite-difference operator or Laplacian operator, so as to get a smoothed model in the inversions. Another way of expressing smoothness can be accomplished by using filter operators, among which is the recursive filter (Lorenc 1992). In general, the application of a one-dimensional first-order recursive filter has two steps (Purser et al., 2003):

$$q_i = (1 - \alpha_i)p_i + \alpha_i q_{i-1} \tag{7}$$

$$s_i = (1 - \alpha_i)q_i + \alpha_i s_{i+1} \tag{8}$$

where  $p$  is the input field,  $s$  is the output field,  $q$  is an intermediate state, and subscript  $i$  indicates the grid index.  $\alpha$  is the smooth factor which lies between 0 and 1. A simple illustration of the recursive filter is given in Fig. 1

A compact matrix form can be derived for the recursive filter which is given by (Purser et al. 2003)

$$\mathbf{A}\mathbf{S} = \mathbf{P} \tag{9}$$

Here  $\mathbf{P}$  is the input field vector, and  $\mathbf{S}$  is the output field vector. From Eq. 6, if we define

$$\mathbf{C}_m^{1/2} = (\mathbf{A}\mathbf{B})^{-1} \tag{10}$$

the inverted model  $\mathbf{m}$  should be smoothed by the recursive filter. In 3D, Eqs. 7 and 8 will be successively used along the horizontal directions for every slice and along the vertical direction for every layer. Moreover, we can apply the first-order recursive filter  $n$  times repeatedly to get a roughly equivalent higher-order one (Purser and Wu et al. 2003). And this number  $n$  is termed as smooth number.

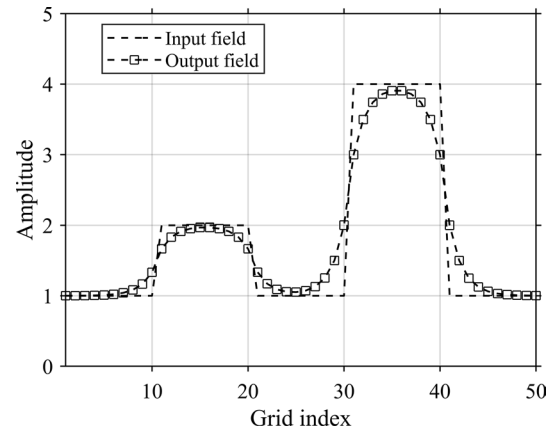


Fig. 1 Illustration of the recursive filter. Apparently, the recursive filter produces a smooth local weighted average of the input field

### Major differences between L-BFGS and NLCG

In this section, we will state the major differences between the L-BFGS and NLCG methods. The first one is the search direction. For NLCG, the search direction is determined by Kelbert et al. (2014)

$$\begin{aligned} \mathbf{h}_k &= -\mathbf{g}_k + \beta \mathbf{h}_{k-1} \\ \beta &= \mathbf{g}_k^T (\mathbf{g}_k - \mathbf{g}_{k-1}) / \mathbf{g}_{k-1}^T \mathbf{g}_{k-1} \end{aligned} \tag{11}$$

where  $\mathbf{h}_k$  and  $\mathbf{h}_{k-1}$  are the current and last model search directions.  $\mathbf{g}_k$  and  $\mathbf{g}_{k-1}$  are the current and previous gradient vectors. For L-BFGS, we require the recent several points and gradients to approximate the inverse Hessian matrix  $\mathbf{H}_k$  and construct the search direction by the following formula (Nocedal and Wright 2006),

$$\mathbf{h}_k = -\mathbf{H}_k \mathbf{g}_k \tag{12}$$

The NLCG update formula is regarded as an extremely L-BFGS method only when giving up information about previous model difference and gradient difference pairs (Koyama et al. 2014).

The second difference is their initial trial step length for every iteration. At the first iteration in both methods, the first trial step length can be computed by

$$\alpha = l / \|\mathbf{g}_0\|_2 \tag{13}$$

where  $\alpha$  is the step length and  $l$  is a constant value. For the rest of iterations, their trial step lengths differ from each other and are given by (Nocedal and Wright 2006)

$$\alpha_{lbfgs} = 1.0$$

$$\alpha_{nlcg} = 1.01 \times \frac{2(f_k - f_{k-1})}{\mathbf{g}_{k-1}^T \mathbf{h}_{k-1}} \quad (14)$$

Different strategies of the initial trial step length are based upon the converge properties of these methods.

The last important difference is the termination of the line search. For a successful line search, the Armijo rule ensures sufficient decrease of the functional value, which is

$$f(\mathbf{m}_k + \alpha_k \mathbf{h}_k) \leq f(\mathbf{m}_k) + c_1 \alpha_k \mathbf{g}_k^T \mathbf{h}_k \quad (15)$$

Here  $c_1$  generally equals 0.0001. Another condition is the so-called curvature condition, expressed as

$$\mathbf{g}(\mathbf{m}_k + \alpha_k \mathbf{h}_k)^T \mathbf{h}_k \leq c_2 \mathbf{g}_k^T \mathbf{h}_k \quad (16)$$

## The L-BFGS inversion framework

While convergence rate is of great concern for 3D inverse problem, we will make two modifications in the conventional L-BFGS algorithm which is given below. For the NLCG algorithm in ModEM, we noticed that the initial trial step length at the first iteration, i.e.,  $\alpha_{0,0}$  in algorithm 1, is scaled by a real scalar (Kelbert et al. 2014). And the search direction at the first iteration is the same in both methods, which is the steepest descent direction. However, in L-BFGS, this entity is fixed in algorithm 1. So, we attempted to change it to the same form as that of NLCG according to Eq. 13.

---

### Algorithm 1 The conventional L-BFGS algorithm

---

Choose starting point  $\mathbf{x}_0$

do  $k = 0, 1, 2 \dots$

    Compute the objective functional value  $f_{ake}$  and gradient  $\mathbf{g}_k$

    Calculate the search direction  $\mathbf{h}_k$  via a two-loop recursion

    Line search: find a step length  $\alpha_k$  where it satisfies conditions 15 and 16

        Initial trial:  $\alpha_{0,0} = 1/\|\mathbf{g}_0\|_2$ , otherwise try  $\alpha_{k,0} = 1.0$  if  $k > 0$

        do  $i = 0, 1, \dots, l_{max}$

            Check on conditions 15 and 16

            Get a new step length  $\alpha_{k,i+1}$  using a line search method.

        End do

    Update  $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{h}_k$

    Check convergence

End do

---

And  $c_2$  generally equals 0.9. For the NLCG method in ModEM, line search procedure just needs to find a step length satisfying condition 15. However, in Nocedal's L-BFGS algorithm, both conditions must be satisfied. From condition 15, in order to check whether a new step length is suitable or not, the objective functional must be evaluate again, which will cost a significant amount time if the forward computation is expensive. Additionally, if we want to find a step length restricted to condition 16, a new gradient needs to be calculated. Therefore, it will be more computationally expensive for the L-BFGS method if an equal number of line search runs at an iteration in both methods, which leads to our proposition of the relaxation line search in this paper.

The other modification relates to the line search procedure. We find that sometimes the line search procedure might not find an acceptable step length that satisfies conditions 15 and 16 quickly. Reducing the maximum line search number, i.e.,  $l_{max}$  in algorithm 1 (initially set to be 20), could limit the line search iterates but risk the success of line search procedure. Instead, we limit the line search under 2 iterates by a relaxation strategy as follows:

**Algorithm 2** The proposed relaxation line search strategyLine search: find a proper step length  $\alpha_k$  at the  $k^{\text{th}}$  iterationSet initial trial step,  $\alpha_{0,0} = 1/\|\mathbf{g}_0\|_2$ , otherwise try  $\alpha_{k,0} = 1.0$  if  $k > 0$ do  $i = 0, 1$ 

If conditions 15 and 16 are satisfied, exit line search

    Get a new step length  $\alpha_{k,i+1}$ 

End do

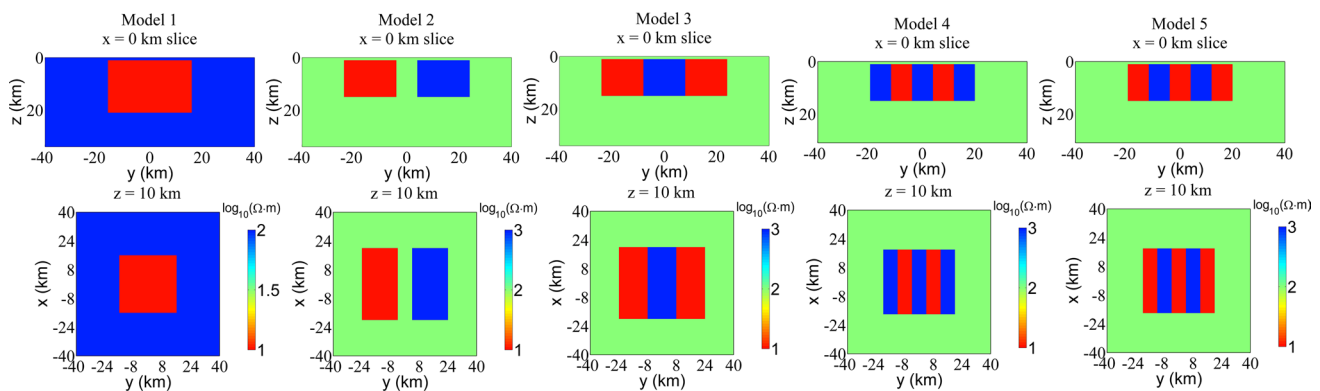
In case the conditions are declined, then select  $\alpha_k$  subject to  $\min(f_{k,0}, f_{k,1})$  where  $f_{k,i} = f(\mathbf{x}_i + \alpha_{k,i})$ 

By using this strategy, two evaluations of the function value and gradient are needed at most. While the selected point will be used at the next iteration of the inversion, the actual calculation requires only one evaluation of the function value and gradient. In such way, the computation workload of the line search is comparable with that implemented in ModEM.

## Results and discussion

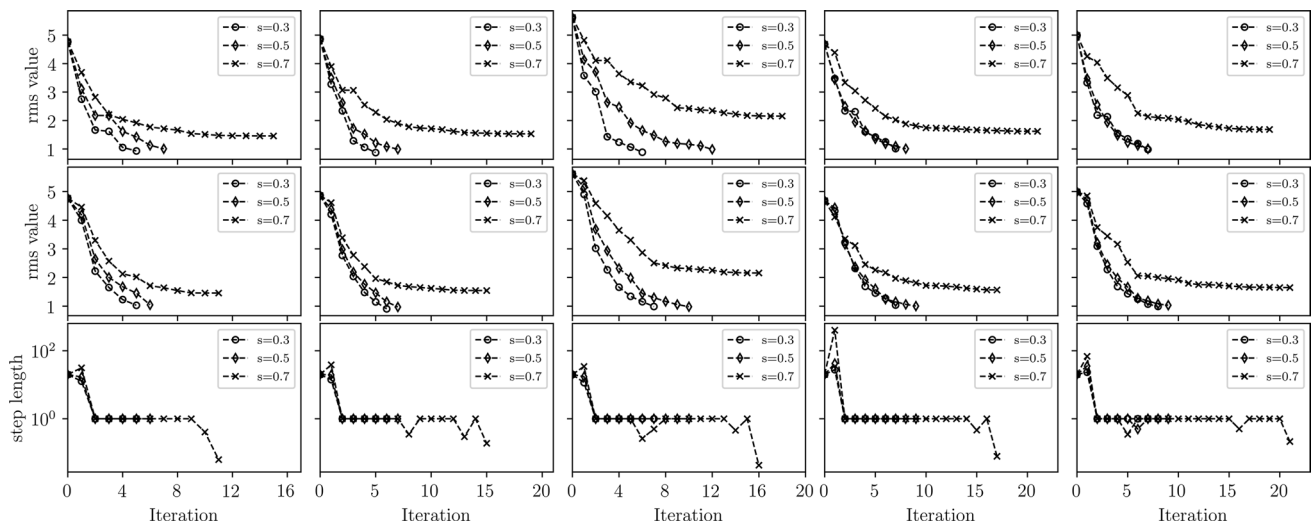
To evaluate our inversion algorithm, we will run a number of synthetic tests on five theoretical models and compare the results with that by the NLCG method in ModEM. And we run these tests under different controlling parameters including the smooth factor and smooth number, the regularization parameter, the initial trial step, and the starting model to investigate the influence of those parameters on the effects of the inversion. The theoretical models are shown in Fig. 2. The first three models are divided into a  $36 \times 36 \times 37$  grid (with 10 air layers). The measured points for those models are located from  $-38$  to  $38$  km along the X direction and

from  $-38$  to  $38$  km along the Y direction, with a distance of 4 km between two adjacent points. For the models 4 and 5, the model domains are divided into  $56 \times 56 \times 41$  cells. And the stations are located from  $-39$  to  $39$  km along the X direction and from  $-39$  to  $39$  km along the Y direction, with a distance of 2 km between each point. The synthetic data are firstly modeled using ModEM which is based on the staggered grid finite difference method. And then 5% of Gaussian noise is added to the impedance tensor components. The periods we selected for all the models are 0.1, 1, 5, 10, 20, 30, 50, 80, 100, and 1000 s. The data variance is set to be 5% of  $|\mathbf{Z}_{xy} \times \mathbf{Z}_{yx}|^{1/2}$ . The grids are the same for the inversion and forward modeling. In addition, the forward iterative solver is terminated when the normalized misfits are below  $10^{-9}$  and  $10^{-8}$  for the first three models and last two models, respectively. And the normalized misfits of the adjoint solver for the gradient calculation are set at  $10^{-6}$  for all models. We determine the inversion should be terminated when the root mean square (RMS) of data misfit is less than 1.05 or cannot further decrease before it reaches target value.



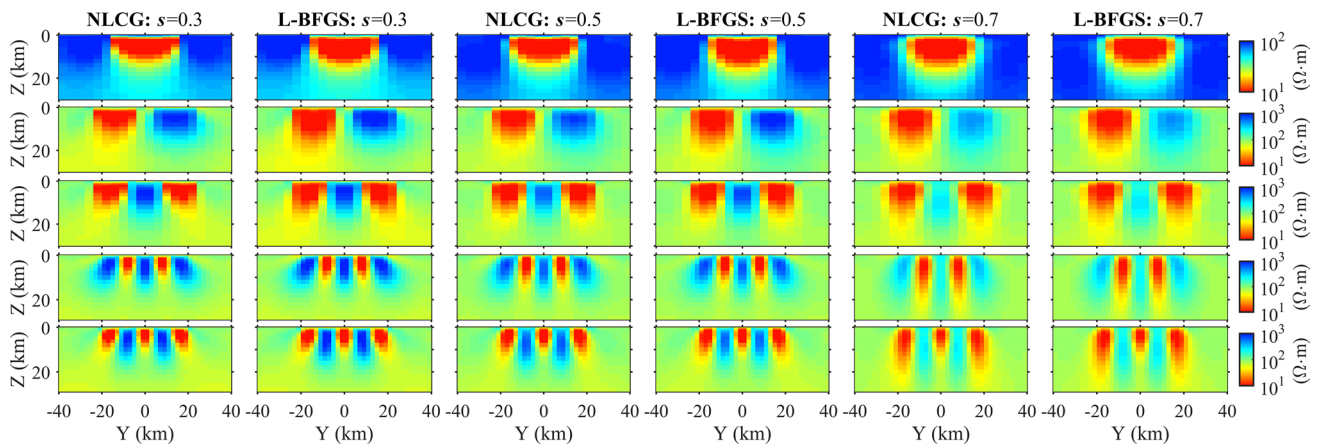
**Fig. 2** The section and plane view of the five theoretical models for the synthetic inversion tests. The background resistivity is  $100 \Omega \text{ m}$ . The low and high resistivities are  $10 \Omega \text{ m}$  and  $1000 \Omega \text{ m}$ , respectively





**Fig. 3** The rms values versus iteration during the synthetic tests for different *smooth factors* in the NLCG inversions (the first row) and L-BFGS inversions (the second row). The third row of panels shows

the step lengths in the L-BFGS inversions. The results for the five models, model 1 to model 5, are displayed sequentially from the left to the right column of panels



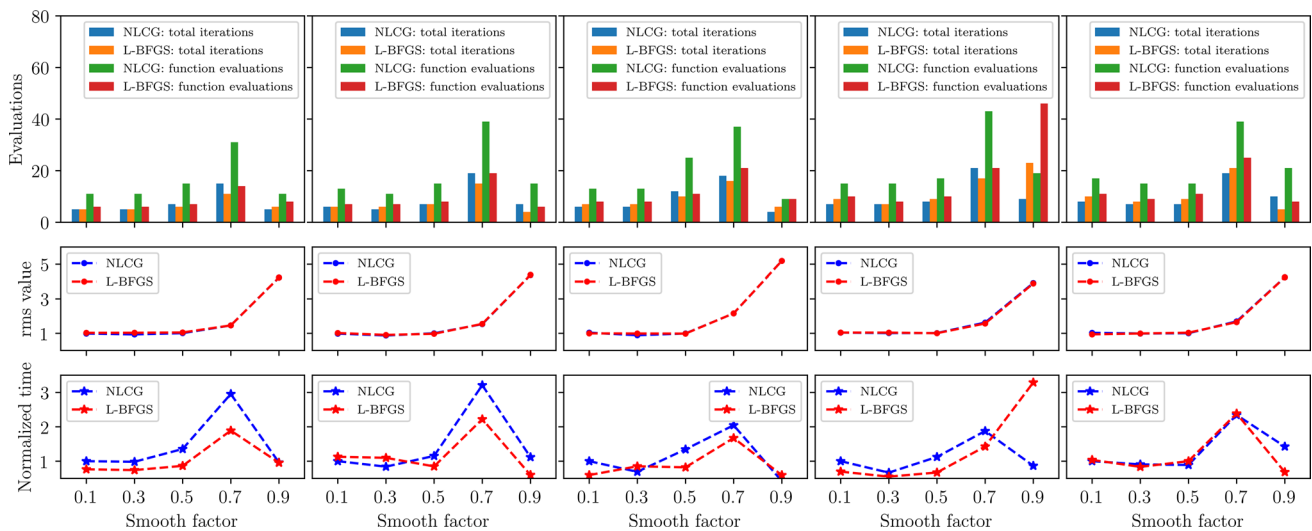
**Fig. 4** The cross-section ( $x=0$  km) of the inverted resistivity models of the synthetic tests for different *smooth factors* in the NLCG and L-BFGS inversions. The results for the five models, model 1 to model 5, are displayed from the top to the bottom of panels

**Synthetic tests for different smooth factors**

As mentioned previously, the model covariance matrix plays an important role in preconditioning the inverse problem. One of the parameters that control the derivation of the matrix is the smooth factor. Therefore, we will run tests for different smooth factors to see how it affects the inversion result. The regularization parameter and the smooth number are fixed at 1.0. The initial trial step length of the first iteration is  $10/\|g_0\|_2$  and the starting model is a  $50 \Omega \cdot m$  half space. The curves of rms values and step lengths are plotted in Fig. 3. The final inverted resistivity models are shown in

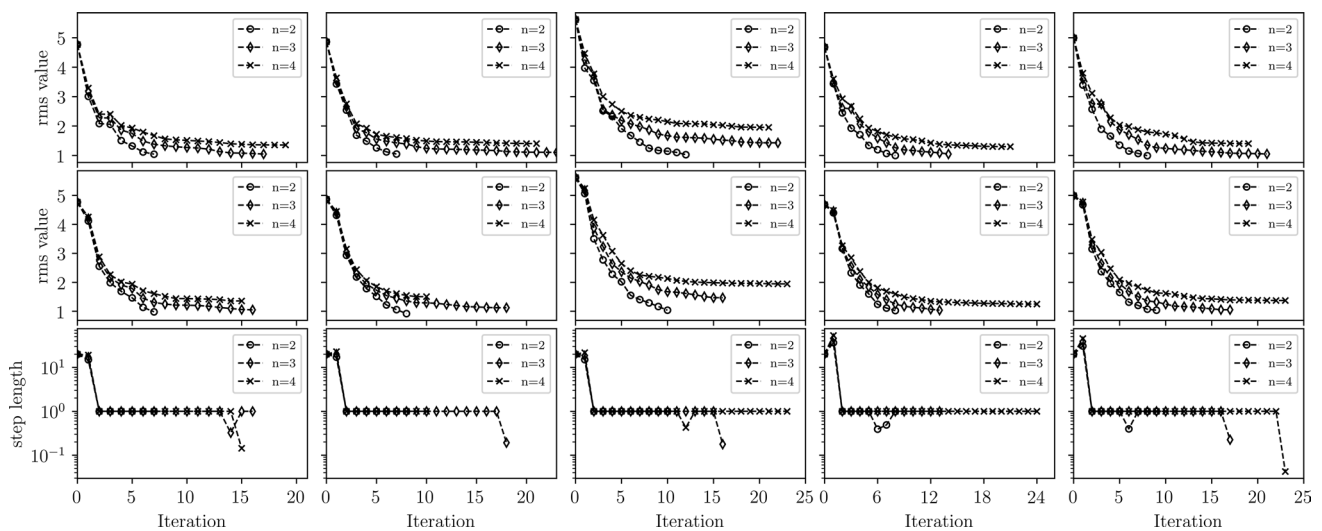
Fig. 4. And the information about computation cost is given in Fig. 5. The smooth factor is notated by  $s$  in all the figures.

Generally, the inversion under larger smooth factor gives us smoother model. However, it suggests that larger smooth factors slow down the convergence rate in Fig. 3. And setting smooth factor greater than 0.5 should be avoided since the rms misfit could not converge to the target value which will result in underfitted models, as shown in Fig. 4. Little differences on the inverted resistivity models are shown by the NLCG and L-BFGS inversions. For the computation cost, the L-BFGS inversions are comparable with the NLCG’s results. In some cases, the L-BFGS inversions turn out to be



**Fig. 5** The panels in the first row show the total iterations and function evaluations in the NLCG and L-BFGS inversions for different *smooth factors*. The corresponding final rms values and computation

time are plotted in the second and third row of panels, respectively. Note that the computation time is normalized by the time used in the NLCG inversion for the first smooth factors, i.e., 0.1 in the examples



**Fig. 6** The rms values versus iteration during the synthetic tests for different *smooth numbers* in the NLCG inversions (the first row) and L-BFGS inversions (the second row). The third row of panels shows

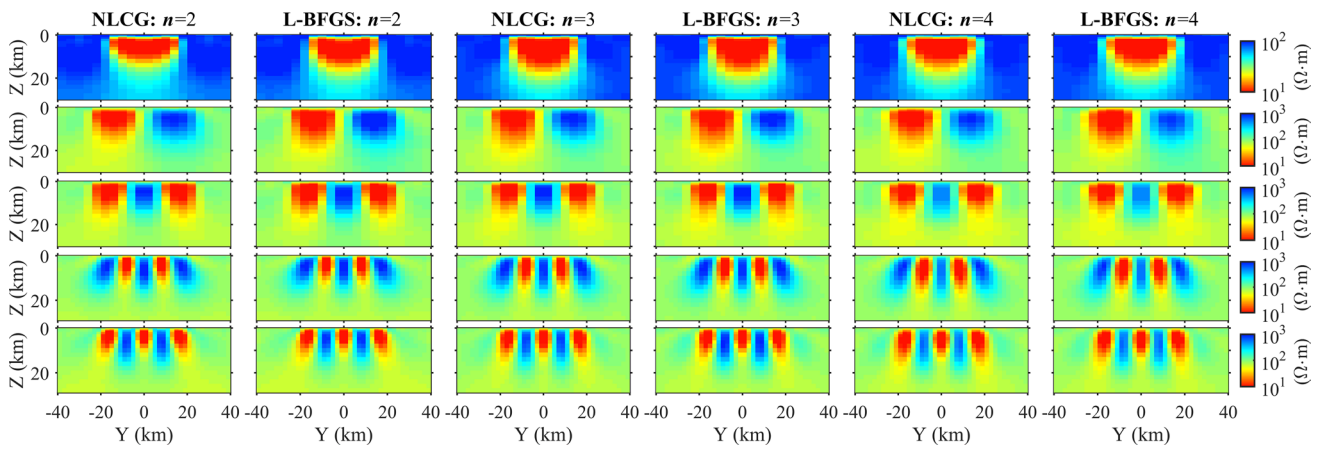
the step lengths in the L-BFGS inversions. The results for the five models, model 1 to model 5, are displayed sequentially from the left to the right column of panels

slightly more efficient which is due to the line search strategies as we infer. For NLCG, the line search procedure needs to find a suitable step length. For L-BFGS, the step length is 1.0, which is just the initial trial, during the majority of the iterations.

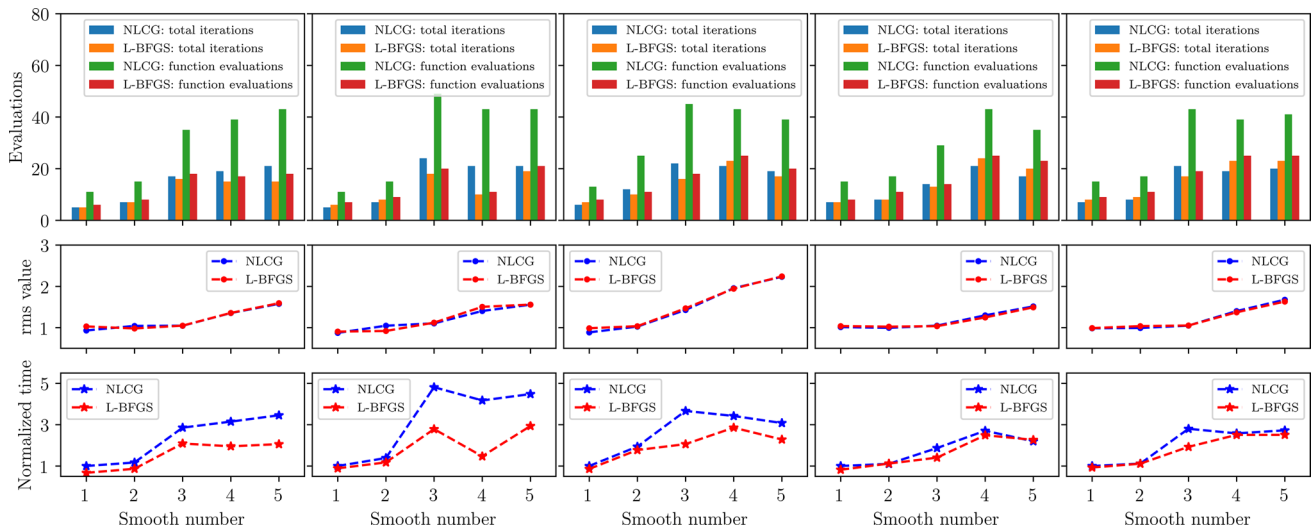
**Synthetic tests for different smooth numbers**

As with the smooth factor, the smooth number described earlier also controls the model covariance. So, the inversion

tests for different smooth numbers are investigated as well. For all the inversion tests, we fix the regularization parameter at 1.0 and the smooth factor at 0.3. The initial trial step length of the first iteration is  $10/\|g_0\|_2$  and the starting model is a 50 Ω m half space. The curves of rms values and step lengths are plotted in Fig. 6. The final inverted resistivity models are shown in Fig. 7. And the information about computation cost is given in Fig. 8. The smooth number is notated by *n* in all the figures.



**Fig. 7** The cross-section ( $x=0$  km) of the inverted resistivity models of the synthetic tests for different *smooth numbers* in the NLCG and L-BFGS inversions. The results for the five models, model 1 to model 5, are displayed from the top to the bottom of panels



**Fig. 8** The panels in the first row show the total iterations and function evaluations in the NLCG and L-BFGS inversions for different *smooth numbers*. The corresponding final rms values and computation time are plotted in the second and third row of panels, respec-

tively. Note that the computation time is normalized by the time used in the NLCG inversion for the first smooth number, i.e., 1 in the examples

From the results, we see that larger smooth number will slow down the convergence rate, as shown in Fig. 6. Setting the smooth number to be 1 or 2 shall be a proper choice since larger values prevent the rms misfit from converging to the target level. The inverted resistivity models by both NLCG and L-BFGS inversions are very similar in Fig. 7. As with the smooth factor, the L-BFGS inversion tends to be slightly more efficient.

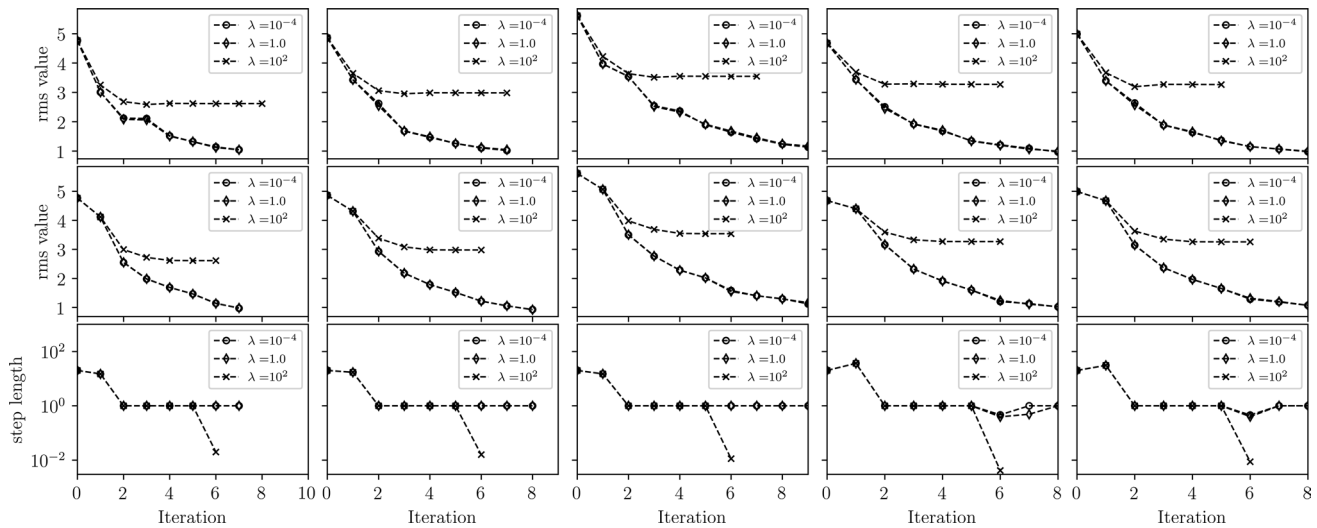
### Synthetic tests for different regularization parameters

As we know, the regularization parameter controls the balance between data term and model regularization term in the objective functional. Large regularization parameter could be used to prevent the inversion from overfitting and solve the inverse problem when it is ill-posed. So, we investigate the impact of different regularization parameters on the inversion. The other controlling parameters are set as

follows: the smooth factor is 0.3, and the smooth number is 2; the initial trial step length of the first iteration is  $10/\|g_0\|_2$ , and the starting model is a 50  $\Omega$  m half space. The curves of rms values and step lengths are plotted in Fig. 9. The final inverted resistivity models are shown in Fig. 10. And the information about the computation cost is given in Fig. 11.

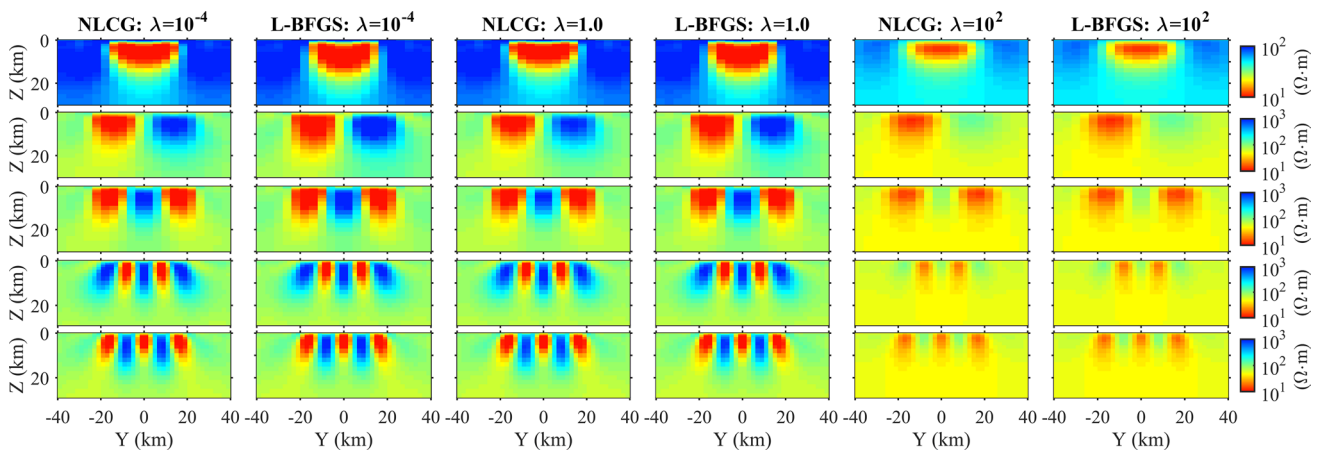
We see that an extremely large regularization parameter will slow down the converge rate in Fig. 9. And the data misfit might fail to reduce to the target level. That is because too much constraint is taken into the model while fitting the data. Additionally, it is surprising that the inverted model is nearly insensitive to small regularization

parameters. We interpret this phenomenon as a benefit from the preconditioning (as in Eq. 5) which produces stable convergence behavior even for the ill-posed problem. And the data misfit is safeguarded above the target level. However, it does not imply that only the small regularization parameters should be used. The selection of the regularization parameters depends on the data and problem. A useful strategy is by using a cooling strategy in which the regularization parameter starts from large values and decreases if the data cannot be fitted properly. As we notice, the recovery of the conductors performs better than that of resistors. And the

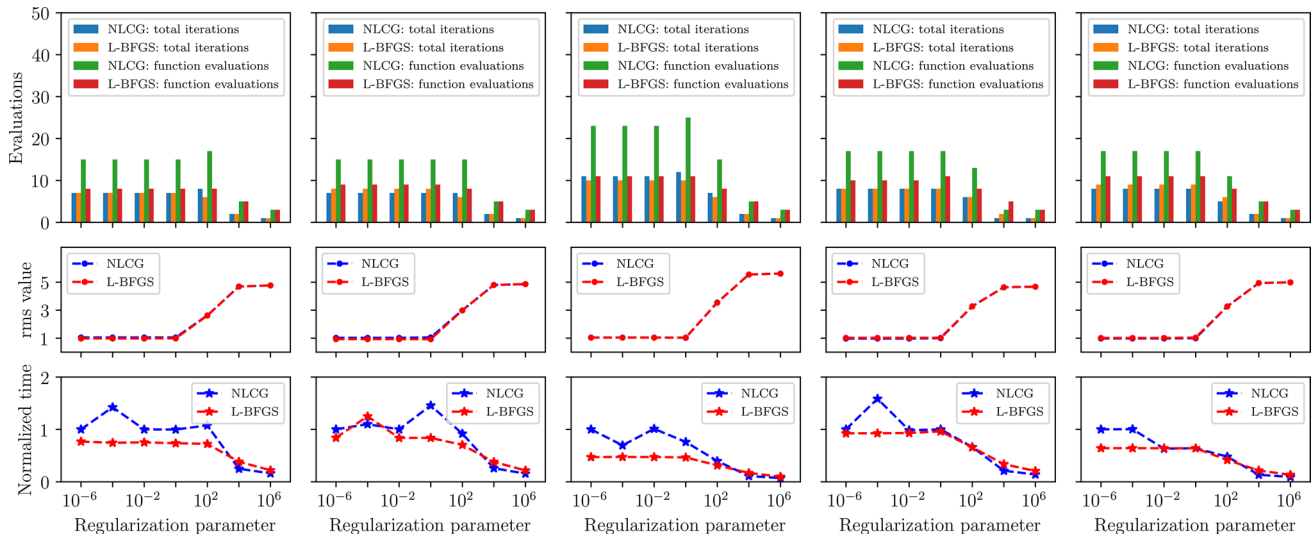


**Fig. 9** The rms values versus iteration during the synthetic tests for different regularization parameters in the NLCG inversions (the first row) and L-BFGS inversions (the second row). The third row of pan-

els shows the step lengths in the L-BFGS inversions. The results for the five models, model 1 to model 5, are displayed sequentially from the left to the right column of panels

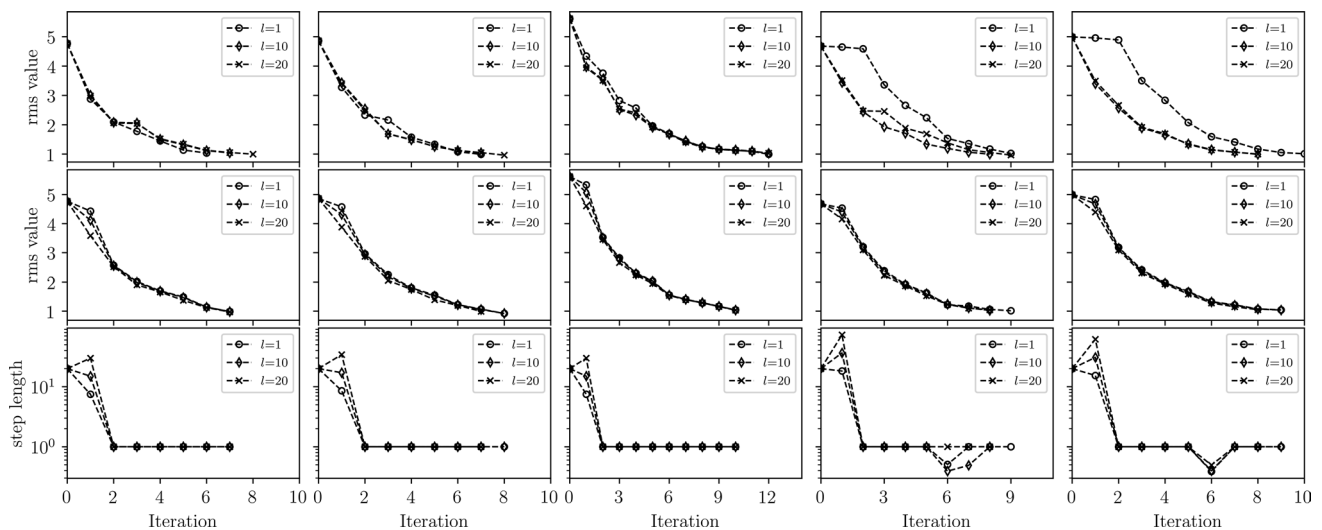


**Fig. 10** The cross section ( $x=0$  km) of the inverted resistivity models of the synthetic tests for different regularization parameters in the NLCG and L-BFGS inversions. The results for the five models, model 1 to model 5, are displayed from the top to the bottom of panels



**Fig. 11** The panels in the first row show the total iterations and function evaluations in the NLCG and L-BFGS inversions for different regularization parameters. The corresponding final rms values and computation time are plotted in the second and third rows of panels,

respectively. Note that the computation time is normalized by the time used in the NLCG inversion for the first regularization parameter, i.e.,  $10^{-6}$  in the examples



**Fig. 12** The rms values versus iteration during the synthetic tests for different initial trial step lengths in the NLCG inversions (the first row) and L-BFGS inversions (the second row). The third row of pan-

els shows the step lengths in the L-BFGS inversions. The results for the five models, model 1 to model 5, are displayed sequentially from the left to the right column of panels

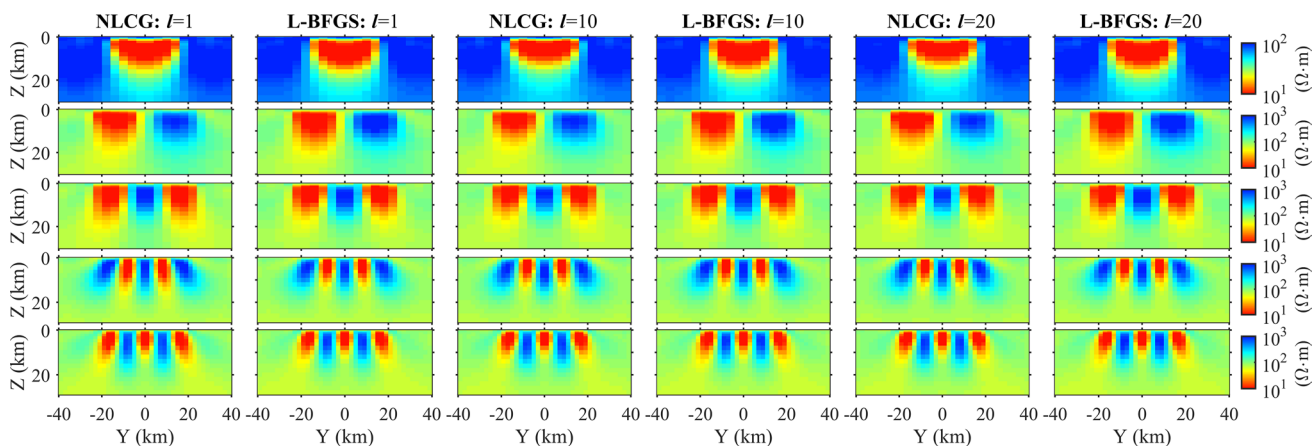
L-BFGS inversion shows considerable improvement in computation time compared with the NLCG inversion.

**Synthetic tests for different initial trial step lengths**

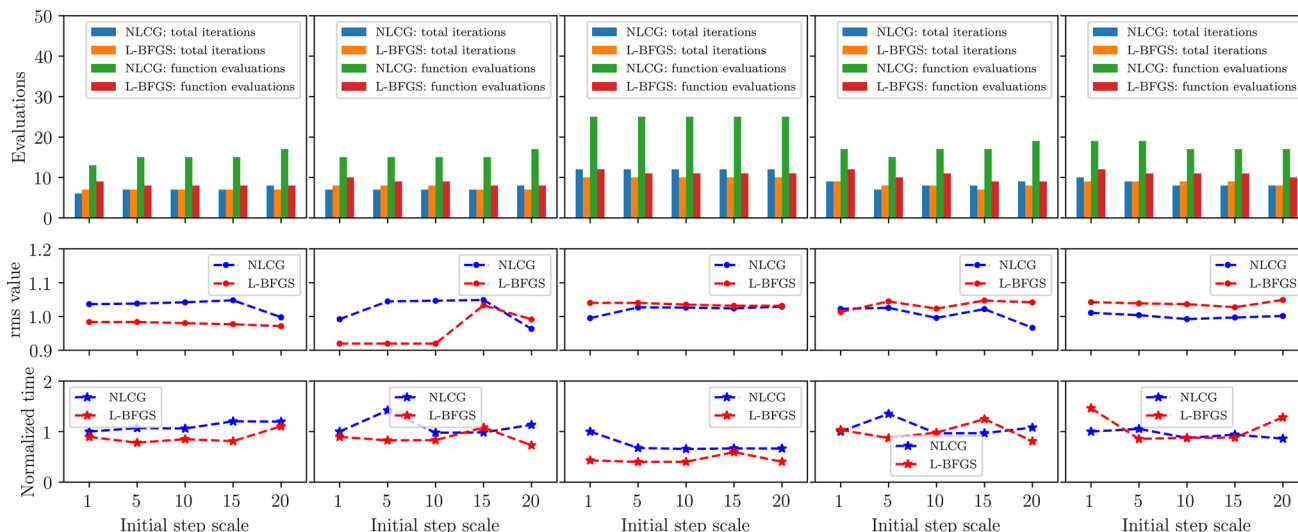
In this section, we test for different initial trial step lengths which have been described previously. The other controlling parameters are set as follows: the smooth factor is 0.3 and the smooth number is 2; the regularization parameter is fixed at 1.0 and the starting model is a 50 Ω m half space.

The curves of rms values and step lengths are plotted in Fig. 12. The final inverted resistivity models are shown in Fig. 13. And the information about the computation cost is given in Fig. 14.

As can be seen from Fig. 13, we get similar inversion results from the L-BFGS method with that from the NLCG inversions. From Fig. 12, we see that more decrease in rms value happens during the early stage of the inversion when a larger initial trial step length is set. The results of model 4 and 5 imply that the NLCG inversion tends to be more



**Fig. 13** The cross section ( $x=0$  km) of the inverted resistivity models of the synthetic tests for different *initial trial step lengths* in the NLCG and L-BFGS inversions. The results for the five models, model 1 to model 5, are displayed from the top to the bottom of panels



**Fig. 14** The panels in the first row show the total iterations and function evaluations in the NLCG and L-BFGS inversions for different *initial trial step lengths*. The corresponding final rms values and computation time are plotted in the second and third row of panels,

respectively. Note that the computation time is normalized by the time used in the NLCG inversion for the first initial trial step length, i.e.,  $1/\|g_0\|_2$  in the examples

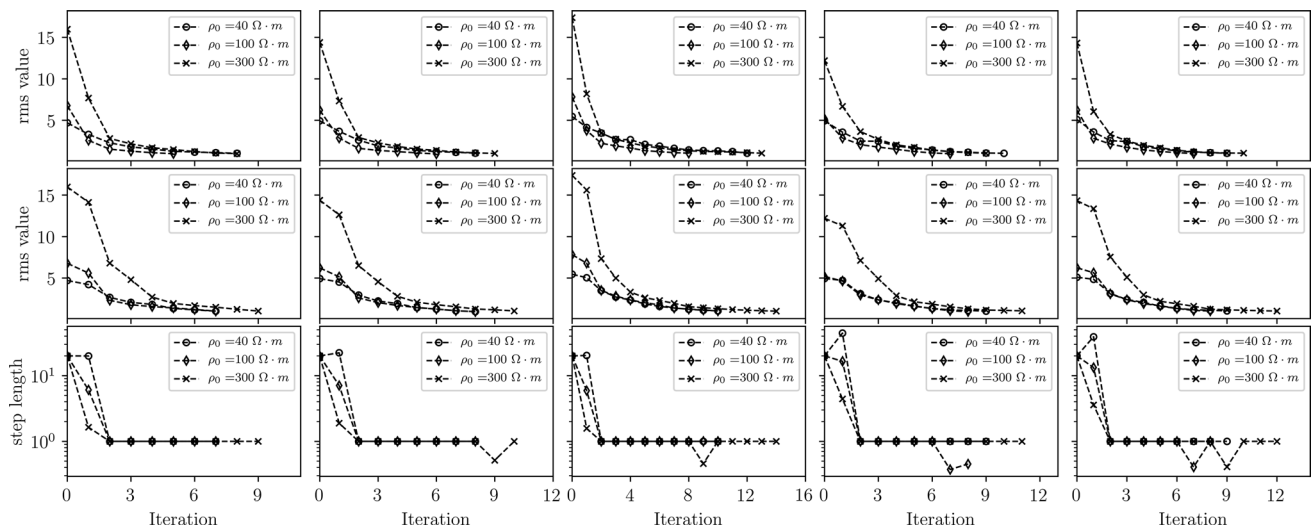
sensitive to this parameter. In most cases, the L-BFGS inversion shows a slight efficiency gain in computation time.

**Synthetic tests for different starting models**

The starting model, which determines the search path of the inversion methods, will be tested for in this section. The other controlling parameters are set as follows: the smooth factor is 0.3 and the smooth number is 2; the regularization parameter is fixed at 1.0. The curves of rms values and step lengths are plotted in Fig. 15. The final inverted resistivity

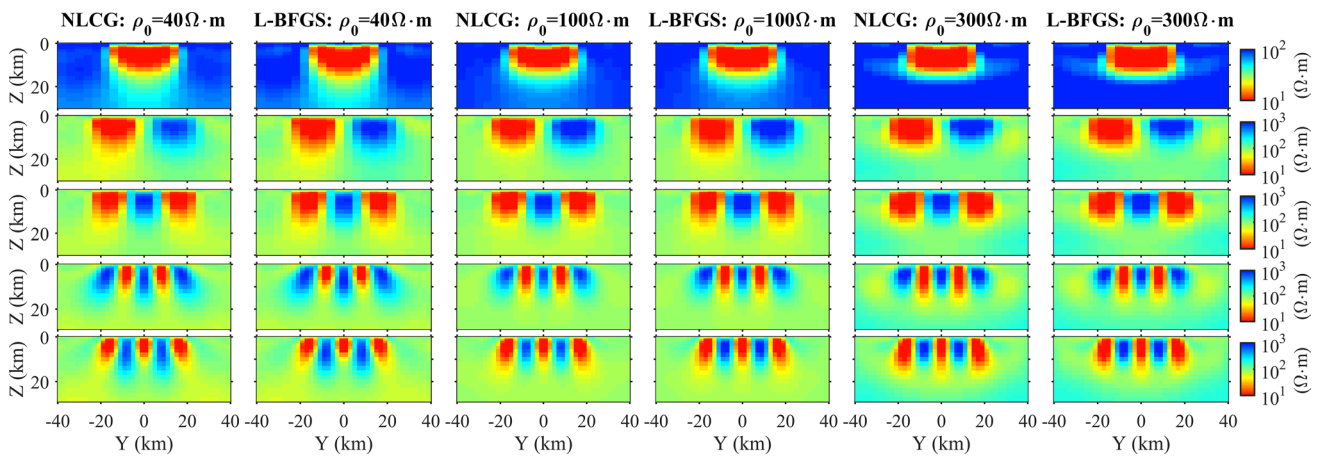
models are shown in Fig. 16. And the information about the computation cost is given in Fig. 17.

In practice, the inversion starting from a model with the background or regional resistivity is preferred. And the rms curves in Fig. 15 show that less iterations are required for such a starting model along with computation time. The inverted models are similar for the L-BFGS and NLCG method, as shown in Fig. 16. Again, the L-BFGS inversion needs less time than NLCG in most cases.



**Fig. 15** The rms values versus iteration during the synthetic tests for different *starting models* in the NLCG inversions (the first row) and L-BFGS inversions (the second row). The third row of panels shows

the step lengths in the L-BFGS inversions. The results for the five models, model 1 to model 5, are displayed sequentially from the left to the right column of panels

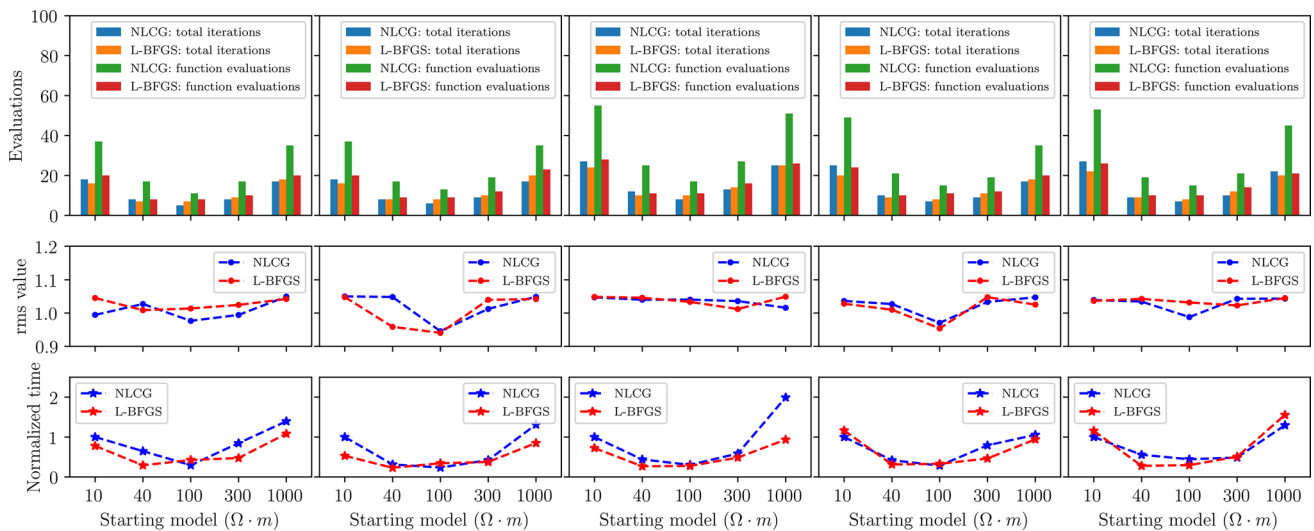


**Fig. 16** The cross section ( $x=0$  km) of the inverted resistivity models of the synthetic tests for different *starting models* in the NLCG and L-BFGS inversions. The results for the five models, model 1 to model 5, are displayed from the top to the bottom of panels

### The real data inversion

To investigate our algorithm on the real data, we perform the L-BFGS inversion on the MT Kayabe dataset (Takasugi et al. 1992). And the inversion with the NLCG method is also conducted as comparison. The dataset consists of 209 sounding points, of which 161 points are highly densely distributed in a rectangular area, as shown in Fig. 18. We choose to invert the data from those highly dense points, excluding 2 irregularly positioned points. Thus, we will invert the data at 159 sites which are equally spaced in the

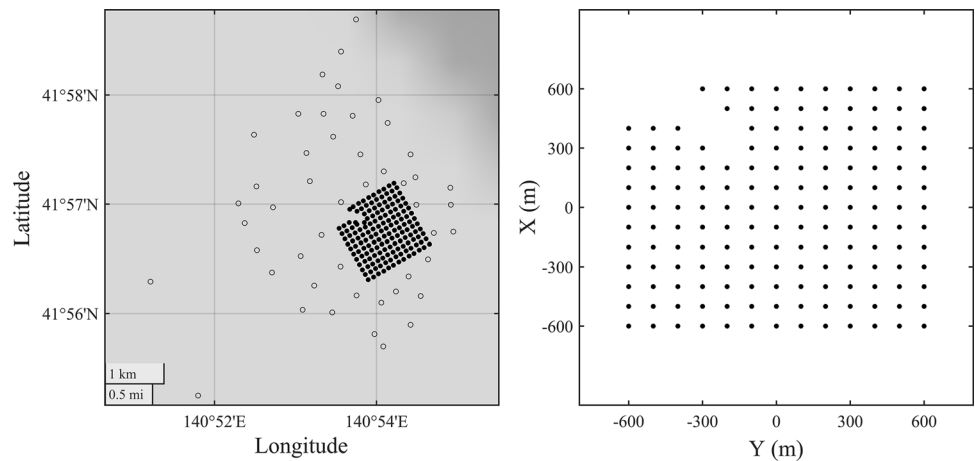
inversion station grid, as shown in Fig. 18. The selected frequencies are ranging within 425 Hz–85.3 s. We divide the model domain into  $35 \times 35 \times 45$  (with 10 air layers) cells. The minimum normalized misfits for the forward solver and adjoint solver are  $10^{-8}$  and  $10^{-6}$ , respectively. The inversion will be terminated when the rms misfit is less than 1.05 or cannot decrease further before reaching the target level. And the error floors are set to be  $0.1 * |\mathbf{Z}_{xy} \times \mathbf{Z}_{yx}|^{1/2}$  for the full impedance tensor components. The maximum number of inversion iterations is set to be 200.



**Fig. 17** The panels in the first row show the total iterations and function evaluations in the NLCG and L-BFGS inversions for different starting models. The corresponding final rms values and computation time are plotted in the second and third row of panels, respectively.

Note that the computation time is normalized by the time used in the NLCG inversion for the first one, i.e., a 10 Ω m half space in the examples

**Fig. 18** The station locations of the Kayabe MT dataset in geographic coordinate system (left panel) and in inversion coordinate system (right panel). The MT sites denoted by solid dots (excluding 2 irregularly-positioned points) are used in this paper, while those in open circles are not



**Table 1** The parameters used in the inversions for the MT Kayabe dataset. P1 and P2 mean two different kinds of parameter settings. The starting models for all the inversions are the half-space with a resistivity of 100 Ω m

Parameter set	Smooth factor	Smooth number	Regularization parameter	Initial trial step length
P1	0.2	2	0.1	$5/\ g_0\ _2$
P2	0.4	3	5	$15/\ g_0\ _2$

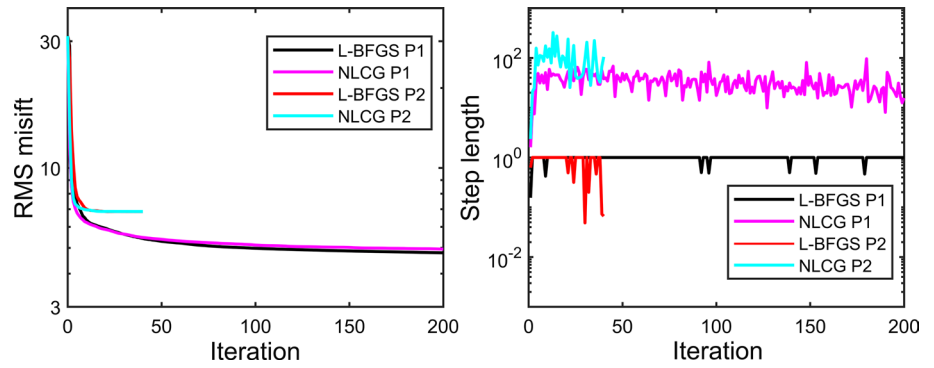
As presented in the previous sections, the inversions under different parameters could lead to inconsistent results. Therefore, we run the inversions for the dataset under two distinct parameter settings, as given in Table 1.

The curves of rms misfit and step lengths during the inversions are plotted in Fig. 19. The statistical information about the inversions is given in Table 2. And the final inverted resistivity models are shown in Fig. 20. From Fig. 19, we see that the inversions under parameter set P2 stops early for both L-BFGS and NLCG, resulting in vague models shown in Fig. 20. Meanwhile, the inversions under parameter P1 reached a lower rms misfit and local heterogeneities are recovered better. More importantly, Table 2 shows that the computation time of the L-BFGS inversions is less than that of NLCG under both parameter settings. Accordingly, the L-BFGS method seems to be more efficient than NLCG for the real data inversion.

As can be seen from Fig. 20, the inversions under both parameter settings show a conductive zone at depth from

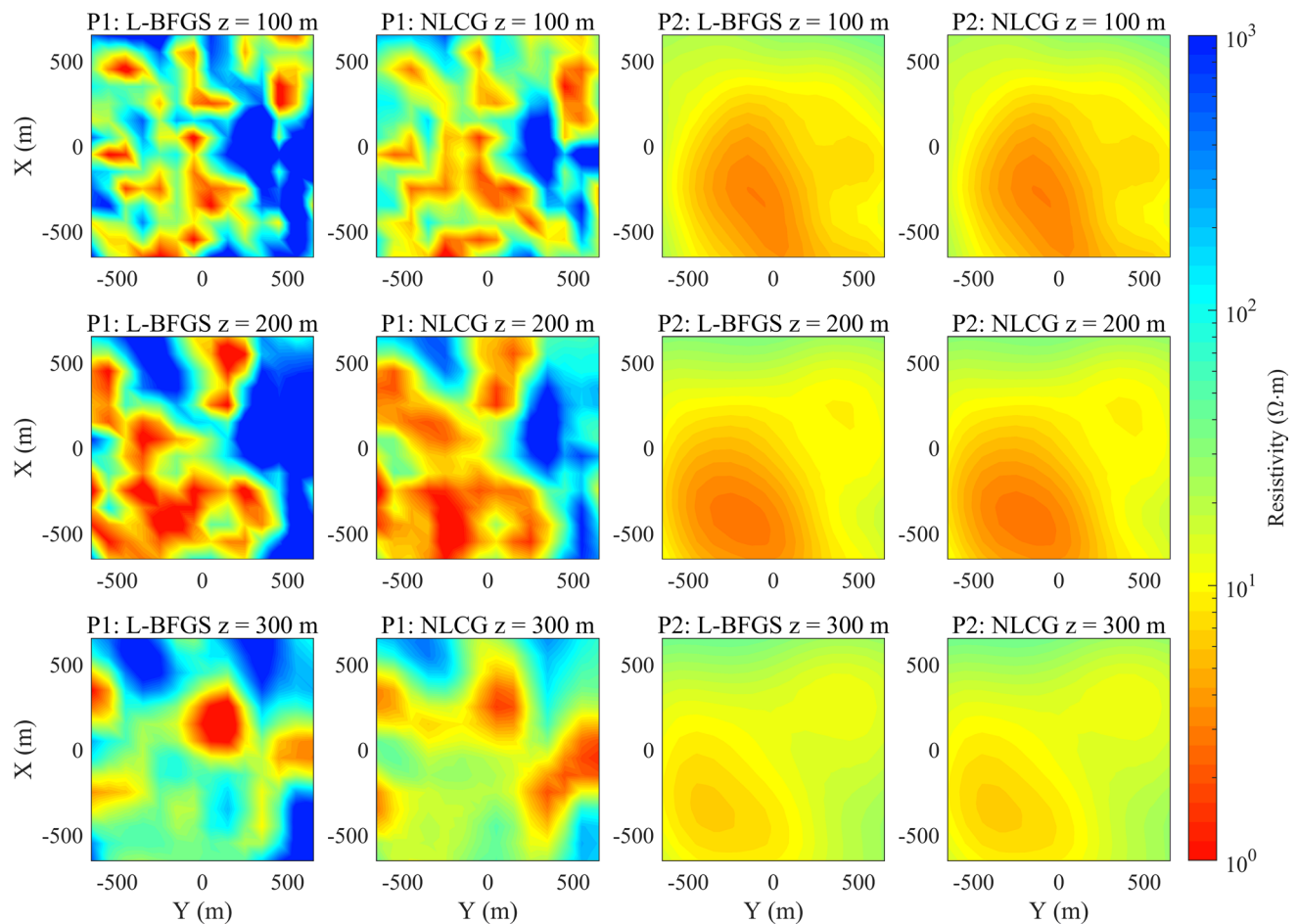


**Fig. 19** The rms curves (left panel) and the step length curves (right panel) for the L-BFGS and NLCG inversions of the Kayabe MT dataset under different controlling parameters. The parameter settings are denoted by P1 and P2

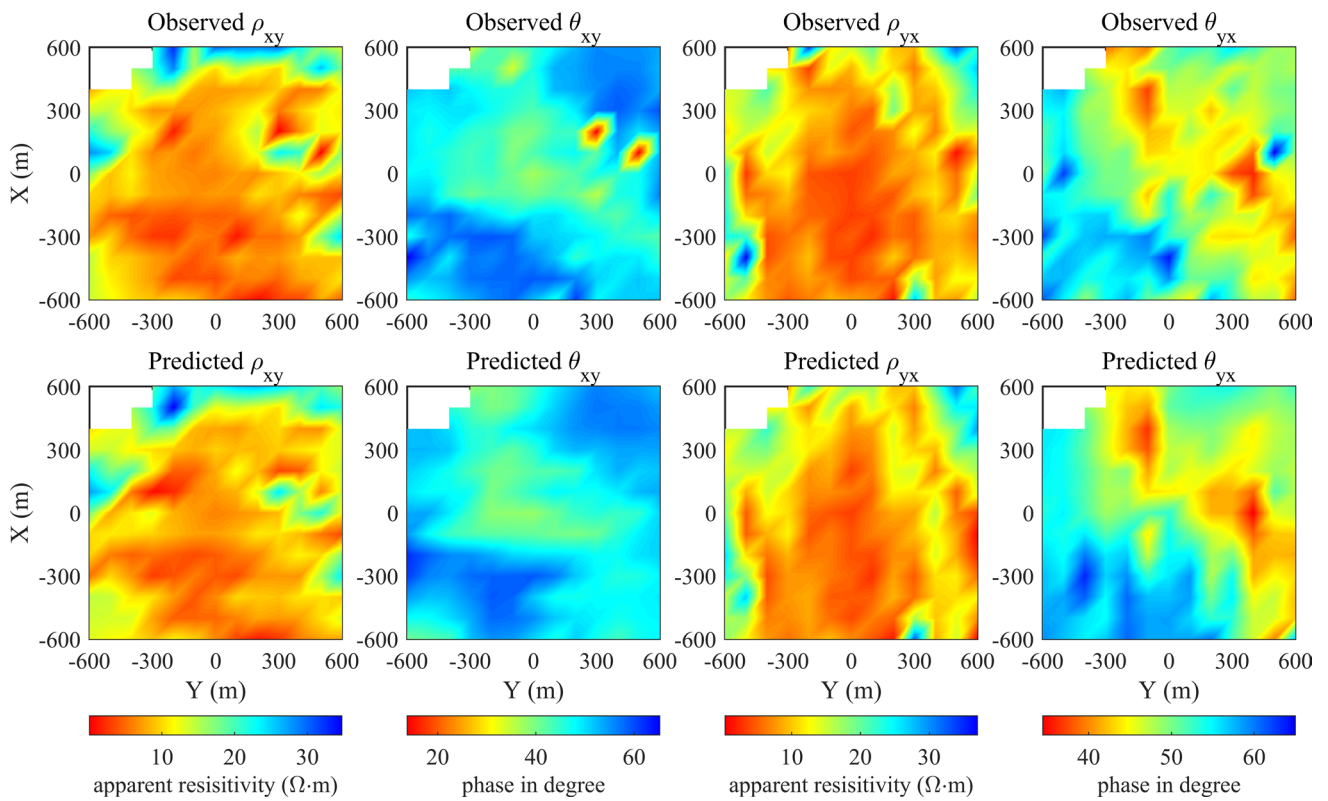


**Table 2** Iteration information and computational cost for the L-BFGS and NLCG inversions under different parameters.  $N_i$ ,  $N_f$ , and  $N_g$  denote total number of iterations, function evaluations, and gradient calculations, respectively

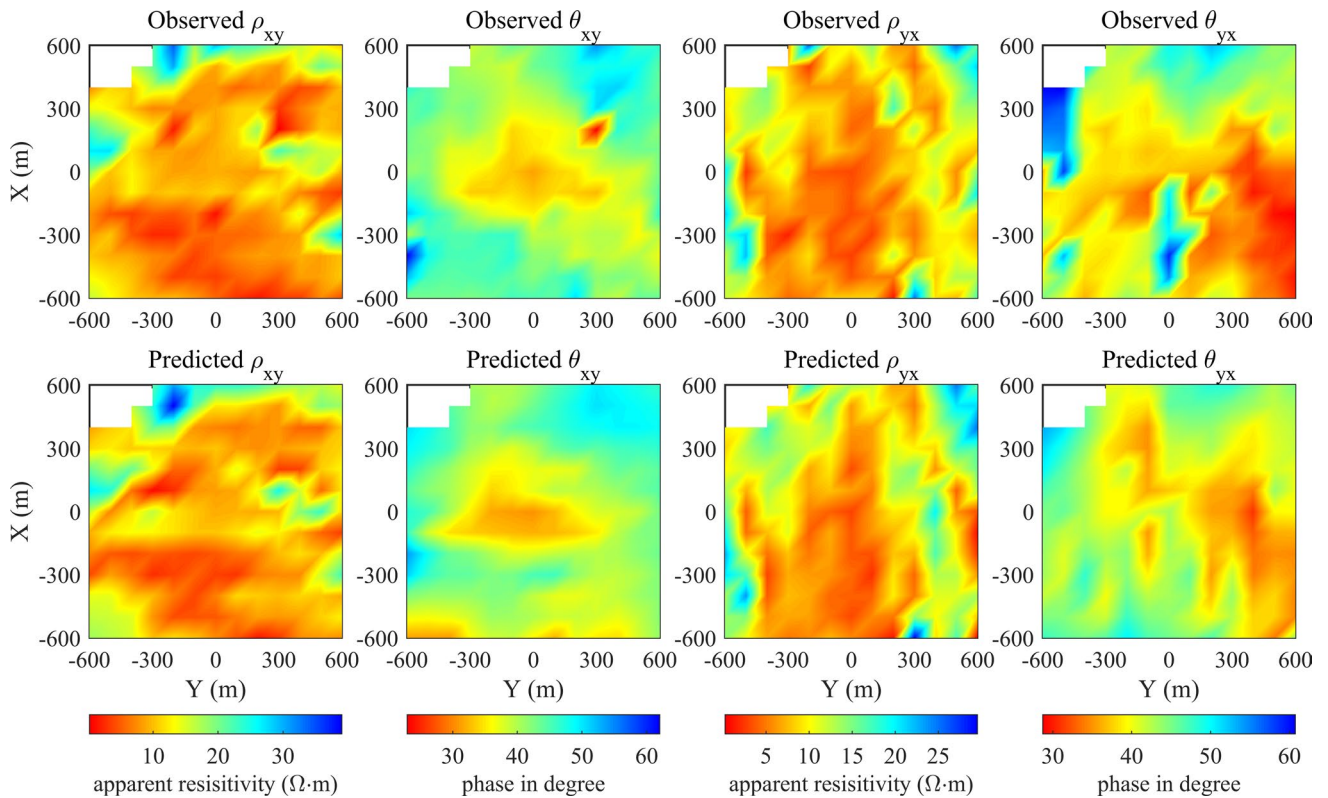
Parameter setting	Method	$N_i$	$N_f$	$N_g$	RMS	Time (h)
P1	L-BFGS	200	207	207	4.79	39.05
	NLCG	200	401	200	4.95	44.17
P2	L-BFGS	40	48	48	6.85	10.78
	NLCG	40	77	40	6.86	11.40



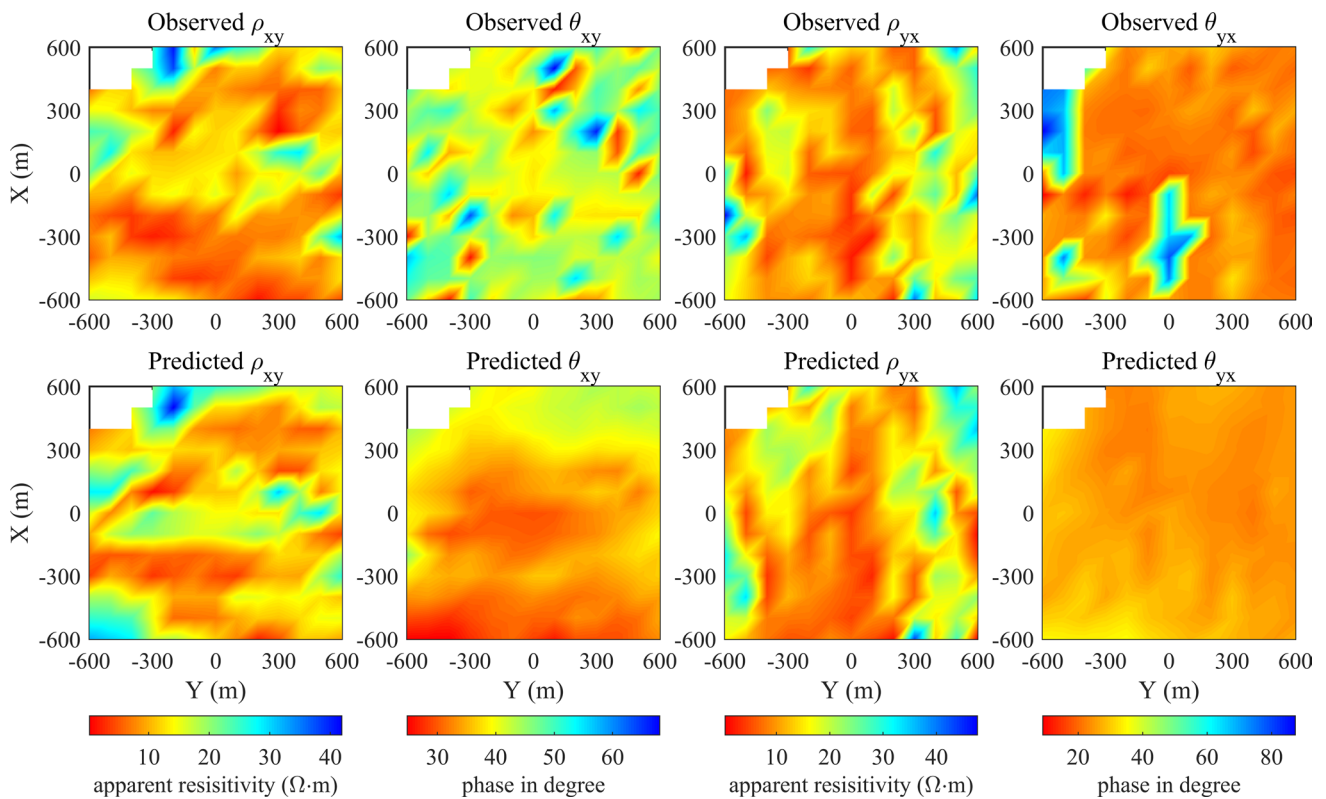
**Fig. 20** The plane views of the final inverted resistivity models for the L-BFGS and NLCG inversions under different parameter settings which are denoted by P1 and P2



**Fig. 21** Comparison of the observed data with the predicted responses of the model inverted by the L-BFGS inversion under the parameter setting P1. The frequency for the data is 24 Hz



**Fig. 22** Comparison of the observed data with the predicted responses of the model inverted by the L-BFGS inversion under the parameter setting P1. The frequency for the data is 8 Hz

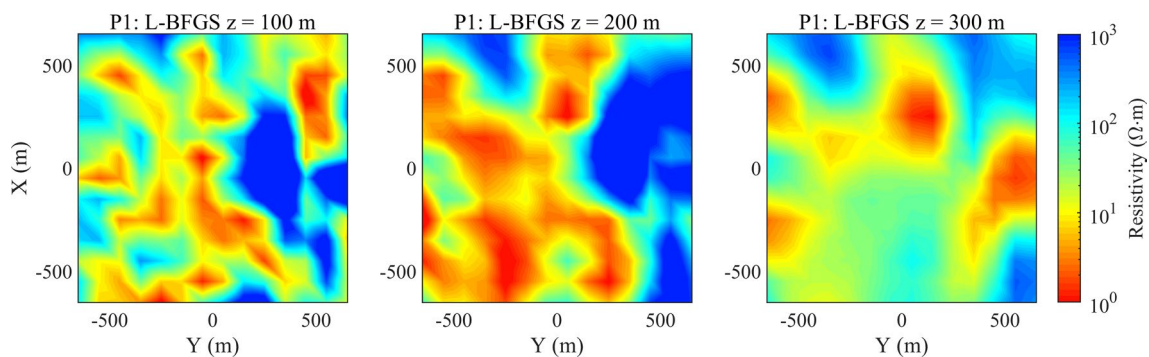


**Fig. 23** Comparison of the observed data with the predicted responses of the model inverted by the L-BFGS inversion under the parameter setting P1. The frequency for the data is 2 Hz

100 to 200 m. The drilling wells indicated the existence of a geothermal reservoir. Lin et al. (2011) inverted the Kayabe data by the conjugate gradient method. To some extent, their results are similar to our inverted model under parameter set P2. We also find the inversion results under parameter P1 show consistency with that by Zhang et al. (2013). To evaluate our inversion results, we compare the observed data with the predicted responses for the inverted model. The comparison results are given in Figs. 21, 22, and 23. As we can see, the data at the frequencies of 24 Hz and 8 Hz are fit considerably well whereas those data at 2 Hz need further improvements. Yamane and Takasugi (1997) suggested the

Kayabe dataset ranging from 3.4 to 250 Hz would provide satisfactory results, which is verified in our inversion. And it might be one of the reasons why the rms misfit could not decrease any further. We are also aware that the data in TM mode are more difficult to fit than that in TE mode, which could be due to the coast effects, as suggested by Takasugi et al. (1992). Therefore, we think our inversion results for the shallow structures should be valid.

Finally, we attempt to interpret the discrepancy of the L-BFGS and NLCG inversions for the real data. From Fig. 19, we can see that the rms misfit in the NLCG inversions is always greater than that in the L-BFGS inversions.



**Fig. 24** The plane views of the inverted resistivity model at the 112th iteration for the L-BFGS inversion under the parameter setting P1

And the inverted models at the 200th iteration for the two methods have some obvious differences. We give the 112th inversion results of L-BFGS inversion in Fig. 24, which is similar to the results at the 200th iteration from NLCG. However, over 80 iterations of the inversion are needed for the NLCG method. We can explain by Eq. 11: when the convergence rate becomes slow and a large number of iterations require during the inversion, the current gradient vector  $\mathbf{g}_k$  and previous gradient vector  $\mathbf{g}_{k-1}$  will gradually become very close, in which situation the search direction of NLCG gradually approximates the minus  $\mathbf{g}_k$ . However, the L-BFGS method uses the last several and current gradients to approximate the inverse Hessian matrix and construct its search directions by Eq. 12. Hence, the L-BFGS method outperforms NLCG under such circumstance due to the second-order sensitivities.

## Conclusions

We have developed a three-dimensional magnetotelluric inversion algorithm based on the L-BFGS method in this paper. The synthetic tests indicate our algorithm is comparable with the NLCG method in ModEM. The recovered resistivity images are similar between one and the other under the same parameter settings. And in most cases, the L-BFGS method could be slightly more efficient than NLCG. The inversion parameters, such as the regularization parameter, should be carefully selected; otherwise, we might fail to get a good interpretation of the subsurface structures. While the ModEM code is widely used among the MT community, our inversion tests for the controlling parameters could be referred to before inverting real data. For the real data inversion, a large number of iterations are needed. Then, not only is the convergence rate of NLCG slow, but also more forward computations are required. However, the L-BFGS method could be more efficient in such situation. Those results show that the L-BFGS algorithm in this paper is a useful and efficient method for 3D MT inverse problem. In addition, our research can be also adopted for other electromagnetic inverse problems whose forward computation is expensive.

**Acknowledgements** This research was funded by the National Natural Science Foundation of China (Grants No. 53200859804 and No. 41830429). The authors are grateful to Nocedal's team for providing the L-BFGS code, Egbert for providing ModEM and NEDO for providing the Kayabe dataset. This paper's research is based on their work.

## Compliance with ethical standards

**Conflict of interest** The authors declare that they have no conflict of interest.

## References

- Avdeev D, Avdeeva A (2009) 3D magnetotelluric inversion using a limited-memory quasi-Newton optimization. *Geophysics* 74(3):F45–F57
- Avdeeva A, Avdeev D (2006) A limited-memory quasi-Newton inversion for 1D magnetotellurics. *Geophysics* 71(5):G191–G196
- Avdeeva A, Avdeev D, Jegen M (2012) Detecting a salt dome overhang with magnetotellurics: 3D inversion methodology and synthetic model studies. *Geophysics* 77(4):E251–E263
- Byrd RH, Lu P, Nocedal J et al (1995) A limited memory algorithm for bound constrained optimization. *SIAM J Sci Comput* 16(5):1190–1208
- Byrd RH, Nocedal J, Schnabel RB (1994) Representations of quasi-Newton matrices and their use in limited memory methods. *Math Program* 63(1–3):129–156
- Devi A, Israil M, Gupta PK et al (2019) Transverse tectonics structures in the Garhwal Himalaya Corridor inferred from 3D inversion of magnetotelluric profile data. *Pure Appl Geophys* 176(11):4921–4940
- Egbert GD, Kelbert A (2012) Computational recipes for electromagnetic inverse problems. *Geophys J Int* 189(1):251–267
- Jahandari H, Farquharson CG (2017) 3-D minimum-structure inversion of magnetotelluric data using the finite-element method and tetrahedral grids. *Geophys J Int* 211(2):1189–1205
- Kelbert A, Egbert GD, Schultz A (2008) Non-linear conjugate gradient inversion for global EM induction: resolution studies. *Geophys J Int* 173(2):365–381
- Kelbert A, Meqbel N, Egbert GD et al (2014) ModEM: a modular system for inversion of electromagnetic geophysical data. *Comput Geosci* 66:40–53
- Koyama T, Khan A, Kuvshinov A (2014) Three-dimensional electrical conductivity structure beneath Australia from inversion of geomagnetic observatory data: evidence for lateral variations in transition-zone temperature, water content and melt. *Geophys J Int* 196(3):1330–1350
- Lin C, Tan H, Tong T (2011) Three-dimensional conjugate gradient inversion of magnetotelluric impedance tensor data. *J Earth Sci* 22(3):386–395
- Liu DC, Nocedal J (1989) On the limited memory BFGS method for large scale optimization. *Math Program* 45(1–3):503–528
- Liu Y, Yin C (2013) 3D inversion for frequency-domain HEM data. *Chin J Geophys Chin Ed* 56(12):4278–4287
- Lorenc A (1992) Iterative analysis using covariance functions and filters. *Q J R Meteorol Soc* 118(505):569–591
- Moorkamp M, Heincke B, Jegen M et al (2011) A framework for 3-D joint inversion of MT, gravity and seismic refraction data. *Geophys J Int* 184(1):477–493
- Newman GA, Alumbaugh DL (2000) Three-dimensional magnetotelluric inversion using non-linear conjugate gradients. *Geophys J Int* 140(2):410–424
- Newman GA, Boggs PT (2004) Solution accelerators for large-scale three-dimensional electromagnetic inverse problems. *Inverse Prob* 20(6):S151–S170
- Newman GA, Gasperikova E, Hoversten GM et al (2008) Three-dimensional magnetotelluric characterization of the Coso geothermal field. *Geothermics* 37(4):369–399
- Ni Q, Yuan YX (1997) A subspace limited memory quasi-Newton algorithm for large-scale nonlinear bound constrained optimization. *Math Comput* 66(220):1509–1520
- Nocedal J, Wright S (2006) *Numerical optimization*. Springer, New York, pp 135–163
- Purser RJ, Wu WS, Parrish DF et al (2003) Numerical aspects of the application of recursive filters to variational statistical analysis.

- Part I: spatially homogeneous and isotropic Gaussian covariances. *Mon Weather Rev* 131(8):1524–1535
- Rodi W, Mackie RL (2001) Nonlinear conjugate gradients algorithm for 2-D magnetotelluric inversion. *Geophysics* 66(1):174–187
- Sass P, Ritter O, Ratschbacher L et al (2014) Resistivity structure underneath the Pamir and Southern Tian Shan. *Geophys J Int* 198(1):564–579
- Siripunvaraporn W, Egbert G (2000) An efficient data-subspace inversion method for 2-D magnetotelluric data. *Geophysics* 65(3):791–803
- Siripunvaraporn W, Sarakorn W (2011) An efficient data space conjugate gradient Occam's method for three-dimensional magnetotelluric inversion. *Geophys J Int* 186(2):567–579
- Siripunvaraporn W, Egbert G, Lenbury Y et al (2005) Three-dimensional magnetotelluric inversion: data-space method. *Phys Earth Planet Int* 150(1–3):3–14
- Takasugi S, Tanaka K, Kawakami N et al (1992) High spatial resolution of the resistivity structure revealed by a dense network MT measurement—a case study in the Minamikayabe Area, Hokkaido Japan. *J Geomagn Geoelectr* 44(4):289–308
- Yamane K, Takasugi S (1997) Data processing procedures for Minamikayabe magnetotelluric soundings. *J Geomagn Geoelectr* 49(11–12):1697–1715
- Zhang K, Dong H, Yan J et al (2013) A NLCG inversion method of magnetotellurics with parallel structure. *Chin J Geophys Chin Ed* 56(11):3922–3931



# Prestack AVO inversion for brittleness index of shale based on BI\_Zoeppritz equation and NSGA II

Chenchen Bi<sup>1</sup> · Yanchun Wang<sup>1</sup> · Wei Xie<sup>2</sup> · Wei Sun<sup>3</sup> · Wei Liu<sup>1</sup>

Received: 12 February 2020 / Accepted: 30 June 2020 / Published online: 9 July 2020  
© Institute of Geophysics, Polish Academy of Sciences & Polish Academy of Sciences 2020

## Abstract

As one of the evaluation characteristics of shale sweet spots, the brittleness index (BI) of shale formations is of great significance in predicting the range of sweet spots, and guiding hydraulic fracturing. Based on the three elastic parameters of P-wave velocity ( $V_p$ ), S-wave velocity ( $V_s$ ) and density obtained by conventional prestack AVO inversion, BI can be calculated indirectly using the Rickman formula. However, the conventional AVO inversion based on Zoeppritz approximation assumes that incident angle is small and elastic parameters change slowly, which affects the inversion accuracy of the three elastic parameters. Additionally, using these three elastic parameters to obtain BI indirectly also leads to cumulative errors of the inversion results. Therefore, we propose an inversion method based on BI\_Zoeppritz equation to directly estimate  $V_p$ ,  $V_s$  and BI. The BI\_Zoeppritz equation is an exact Zoeppritz equation for BI, which is used as the forward operator for the proposed method. The multi-objective function of the inversion method is optimized by a fast nondominated sorting genetic algorithm (NSGA II). An initial model and an optimized search window are used to improve the inversion accuracy. The test results of model data and actual data reveal that this method can directly obtain the BI with high precision. In addition, the stability and noise immunity of the proposed method are verified by the seismic data with random noise.

**Keywords** Brittleness index · Shale · AVO inversion · BI\_Zoeppritz equation · NSGA II

## Introduction

Recently, the exploration and development of unconventional shale reservoirs has attracted great attention from the industry (Glorioso and Rattia 2012; McGlade et al. 2013). Unlike conventional reservoirs based on trap evaluation, the goal of shale oil and gas exploration is to predict sweet spots and define reservoir boundary (Jia 2017). As one of the six characteristics for evaluation of shale sweet spots, the BI of shale formations is significant to predict the distribution scope of sweet spots and guide the design of the hydraulic fracturing operation (Zou et al. 2014).

The shale formations of most shale oil and gas fields in the world contain relatively less clay minerals and more brittle minerals (Jarvie et al. 2007; Fu et al. 2015; Gholami et al. 2016; Pei et al. 2016). The brittle minerals represented by siliceous minerals are easy to break and produce fractures under external force (Zhang et al. 2016). These fractures provide the beneficial channels for gas migration and production. Furthermore, in the process of shale reservoir reconstruction, fracture initiation and propagation decided by formation brittleness are the core problems of hydraulic fracturing (Guo and Zhang 2014). Therefore, it is necessary to find quantitative methods to predict BI of shale reservoirs.

Considering strength, hardness and stress–strain characteristics of shale, BI estimation methods based on mineral composition and Rickman formula are proposed (Jarvie et al. 2007; Rickman et al. 2008; Wang and Gale 2009; Jin et al. 2014a, b). The mineral composition method evaluates BI by the content of brittle minerals. Based on the statistical analysis of Barnett shale in North America, Grigg (2004) concluded that high yield shale has the characteristics of high Young's modulus (YM) and low Poisson's ratio (PR). YM indicates the ability of a shale to maintain fractures

✉ Yanchun Wang  
wyc0426@126.com

<sup>1</sup> School of Geophysics and Information Technology, China University of Geosciences, Beijing 100083, China

<sup>2</sup> Sinopec Geophysical Research Institute, Nanjing 210014, China

<sup>3</sup> Sinopec Exploration and Production Research Institute, Beijing 100083, China

after it is fractured, and PR reflects the ability of shale to rupture under a certain pressure. According to the analysis results of mechanical properties of Barnett shale, Rickman et al. (2008) proposed a formula based on YM and PR to calculate BI. Therefore, the three elastic parameters obtained by conventional AVO inversion can be used to calculate YM and PR, and then BI can be given according to the Rickman formula. Estimation method based on Rickman formula has been the most widely used method to predict shale BI.

Considering the cumulative errors caused by indirect calculation for YM and PR, Zong et al. (2013) derived the YPD approximation from Aki–Richards approximation (Aki and Richards 1980). The YPD approximation can be used to directly estimate YM, PR and density based on Bayesian Frame and AVO inversion method (Yin et al. 2015). However, the AVO inversion based on Zoeppritz approximation assumes that incident angle is small and elastic parameters change slowly (Fang et al. 2016; Xie et al. 2019). In this study, we propose a novel AVO inversion method to directly estimate  $V_p$ ,  $V_s$  and BI. To avoid the limitation of approximations, we use an exact Zoeppritz equation for BI (BI\_Zoeppritz equation) as the forward operator of the novel method. AVO inversion is a kind of nonlinear optimization problem with non-unique solutions (Liu and Wang 2018; Yuan et al. 2019). Therefore, we measure the quality of the inversion results from two aspects: correlation coefficient and root mean square error, so as to establish a multi-objective function. To make the inversion results accord with the real situation as much as possible, we use NSGA II to simultaneously minimize the multi-objective function. To reduce the amount of calculations and improve the accuracy of inversion results, we design an initial model and an optimized search window to constrain the entire inversion process. Model data and actual data are used to test the validity and practicability of the novel method.

## Methodology

### BI\_Zoeppritz equation

As the core of AVO analysis and prestack seismic inversion, the exact Zoeppritz equation (Zoeppritz 1919) can calculate the reflection and transmission coefficient of different incident angles when the three elastic parameters of upper and lower layers are given. To reduce the number of parameters, Fang et al. (2016) expressed the Zoeppritz equation as a function of P-wave reflectivity, S-wave reflectivity, density reflectivity and  $V_p$ – $V_s$  ratio of upper layer. However, it is

necessary to assume that the  $V_p$ – $V_s$  ratio of upper layer is background value. To calculate the reflection coefficient directly from BI, Zoeppritz equation can be represented as a function of  $V_p$ ,  $V_s$  and BI by using the conversion relationship between these elastic parameters.

Rickman formula defines the BI as Eq. 1:

$$BI = \frac{BI_{YM} + BI_{PR}}{2} \quad (1)$$

where  $BI_{YM}$  and  $BI_{PR}$  represent the YM and PR of rock in percentage after normalization, respectively, and the value range is from 0 to 100. This is mainly due to the large dimensional difference between YM and PR, and the data processing of YM and PR is made to reflect their role in the evaluation of BI, as shown in Eqs. 2 and 3:

$$BI_{YM} = 100 \times \frac{YM - YM_{\min}}{YM_{\max} - YM_{\min}}, \quad (2)$$

$$BI_{PR} = 100 \times \frac{PR - PR_{\max}}{PR_{\min} - PR_{\max}}, \quad (3)$$

where  $YM_{\min}$  and  $YM_{\max}$  denote the minimum and maximum of YM; and  $PR_{\min}$  and  $PR_{\max}$  denote the minimum and maximum of PR, respectively. Based on  $V_p$ ,  $V_s$  and density, the YM and PR can be expressed as Eqs. 4 and 5, respectively:

$$YM = 2V_s^2 \rho \left( \frac{3V_p^2 - 4V_s^2}{2V_p^2 - 2V_s^2} \right), \quad (4)$$

$$PR = \frac{V_p^2 + 2V_s^2}{2V_p^2 - 2V_s^2}. \quad (5)$$

where  $\rho$  denotes density. The relationship between BI,  $V_p$ ,  $V_s$  and  $\rho$  can be expressed as Eq. 6:

$$BI = 50 \times \frac{2V_s^2 \rho \left( \frac{3V_p^2 - 4V_s^2}{2V_p^2 - 2V_s^2} \right) - YM_{\min}}{YM_{\max} - YM_{\min}} + 50 \times \frac{\frac{V_p^2 + 2V_s^2}{2V_p^2 - 2V_s^2} - PR_{\max}}{PR_{\min} - PR_{\max}}. \quad (6)$$

Since  $YM_{\min}$ ,  $YM_{\max}$ ,  $PR_{\min}$  and  $PR_{\max}$  are known parameters obtained from prior information of the actual work area, BI can be regarded as the nonlinear function of  $V_p$  and  $V_s$ , and the linear function of  $\rho$ . Similarly,  $\rho$  can also be seen as a function of the BI,  $V_p$  and  $V_s$ , as shown in Eq. 7, which is abbreviated as Eq. 8:

$$\rho = \frac{\left( \frac{BI}{50} - \frac{V_p^2 + 2V_s^2}{2V_p^2 - 2V_s^2} - PR_{\min} - PR_{\max} \right) (YM_{\max} - YM_{\min}) + YM_{\min}}{V_s^2 \left( \frac{3V_p^2 - 4V_s^2}{V_p^2 - V_s^2} \right)}, \tag{7}$$

$$\rho = \phi(BI, V_p, V_s). \tag{8}$$

By replacing all the density items in the exact Zoeppritz equation with Eq. 7, the Zoeppritz equation of  $V_p$ ,  $V_s$  and BI (BI\_Zoeppritz) can be expressed as shown in Eq. 9:

$$\begin{cases} E = b \frac{\cos \alpha_1}{V_{p1}} + c \frac{\cos \alpha_2}{V_{p2}} \\ F = b \frac{\cos \beta_1}{V_{s1}} + c \frac{\cos \beta_2}{V_{s2}} \\ G = a - d \frac{\cos \alpha_1}{V_{p1}} \frac{\cos \beta_2}{V_{s2}} \\ H = a - d \frac{\cos \alpha_2}{V_{p2}} \frac{\cos \beta_1}{V_{s1}} \\ D = EF + GHp^2 \end{cases}, \tag{12}$$

and  $p = \sin \alpha_1 / V_{p1}$  is ray parameter.

We use the Ostrander (1984) shale model listed in Table 1

$$\begin{bmatrix} \sin \alpha_1 & \cos \beta_1 & -\sin \alpha_2 & \cos \beta_2 \\ \cos \alpha_1 & -\sin \beta_1 & \cos \alpha_2 & \sin \beta_2 \\ \sin 2\alpha_1 \frac{V_{p1}}{V_{s1}} \cos 2\beta_1 - \frac{V_{p1} V_{s2}^2}{V_{p2} V_{s1}^2} \phi(BI_2, V_{p2}, V_{s2}) \sin 2\alpha_2 - \frac{V_{p1} V_{s2}}{V_{s1}^2} \phi(BI_2, V_{p2}, V_{s2}) \cos 2\beta_2 \\ \cos 2\beta_1 - \frac{V_{p1}}{V_{s1}} \sin 2\beta_1 - \frac{V_{p2}}{V_{p1}} \phi(BI_2, V_{p2}, V_{s2}) \cos 2\beta_2 - \frac{V_{s2}}{V_{p1}} \phi(BI_2, V_{p2}, V_{s2}) \sin 2\beta_2 \end{bmatrix} \times \begin{bmatrix} R_{pp} \\ R_{ps} \\ T_{pp} \\ T_{ps} \end{bmatrix} = \begin{bmatrix} -\sin \alpha_1 \\ \cos \alpha_1 \\ \sin 2\alpha_1 \\ -\cos 2\beta_1 \end{bmatrix}, \tag{9}$$

where  $V_{p1}, V_{p2}, V_{s1}, V_{s2}, \rho_1$  and  $\rho_2$  denote the  $V_p, V_s$  and density of the upper and lower layers,  $\alpha_1, \beta_1, R_{pp}$  and  $R_{ps}$  represent reflection angles and reflection coefficients of the PP- and PS-waves,  $\alpha_2, \beta_2, T_{pp}$  and  $T_{ps}$  represent transmission angles and transmission coefficients of the PP- and PS-waves, respectively. The relationship between the reflection coefficient and BI can be directly established by using the BI\_Zoeppritz equation. The  $R_{pp}$  can be obtained by solving the Zoeppritz equation. Aki and Richards (1980) gave the analytic expression of  $R_{pp}$  about  $V_p, V_s$  and  $\rho$  (Li et al. 2016). By replacing all the density items in the  $R_{pp}$ , the analytic expression of  $R_{pp}$  about  $V_p, V_s$  and BI can be expressed as shown in Eq. 10:

$$R_{pp} = \left[ \left( b \frac{\cos \alpha_1}{V_{p1}} - c \frac{\cos \alpha_2}{V_{p2}} \right) F - \left( a + d \frac{\cos \alpha_1}{V_{p1}} \frac{\cos \beta_2}{V_{s2}} \right) Hp^2 \right] / D, \tag{10}$$

where

$$\begin{cases} a = (1 - 2V_{s2}^2 p^2) \phi(BI_2, V_{p2}, V_{s2}) - (1 - 2V_{s1}^2 p^2) \phi(BI_1, V_{p1}, V_{s1}) \\ b = (1 - 2V_{s2}^2 p^2) \phi(BI_2, V_{p2}, V_{s2}) + 2V_{s1}^2 p^2 \phi(BI_1, V_{p1}, V_{s1}), \\ c = (1 - 2V_{s1}^2 p^2) \phi(BI_1, V_{p1}, V_{s1}) + 2V_{s2}^2 p^2 \phi(BI_2, V_{p2}, V_{s2}) \\ d = 2(V_{s2}^2 \phi(BI_2, V_{p2}, V_{s2}) - V_{s1}^2 \phi(BI_1, V_{p1}, V_{s1})) \end{cases}, \tag{11}$$

to test the accuracy of BI\_Zoeppritz equation. YM and PR of Ostrander shale model can be calculated by using Eqs. 4 and 5. Assuming that  $YM_{\min} = 1.0 \times 10^{10} \text{ kg m}^{-1} \text{ s}^{-2}$ ,  $YM_{\max} = 1.5 \times 10^{10} \text{ kg m}^{-1} \text{ s}^{-2}$ ,  $PR_{\min} = 0.05$  and  $PR_{\max} = 0.5$ , we can obtain the BI of Ostrander shale model according to Eq. 6.

Figure 1 shows PP-wave reflection coefficients ( $R_{pp}$ ) calculated by the exact Zoeppritz equation, BI\_Zoeppritz equation and Aki–Richards approximation. It is evident that the  $R_{pp}$  of BI\_Zoeppritz equation and Zoeppritz equation is coincided with each other, while the  $R_{pp}$  of Aki–Richards approximation gradually deviates from the exact values of Zoeppritz equation with angle increasing. In addition, Aki–Richards approximation cannot reflect the total reflection phenomenon at  $55^\circ$  as shown in Fig. 1a. The BI\_Zoeppritz equation does not introduce any hypothesis and avoids the loss of accuracy, so we can use it as the forward operator of the AVO inversion to estimate BI.

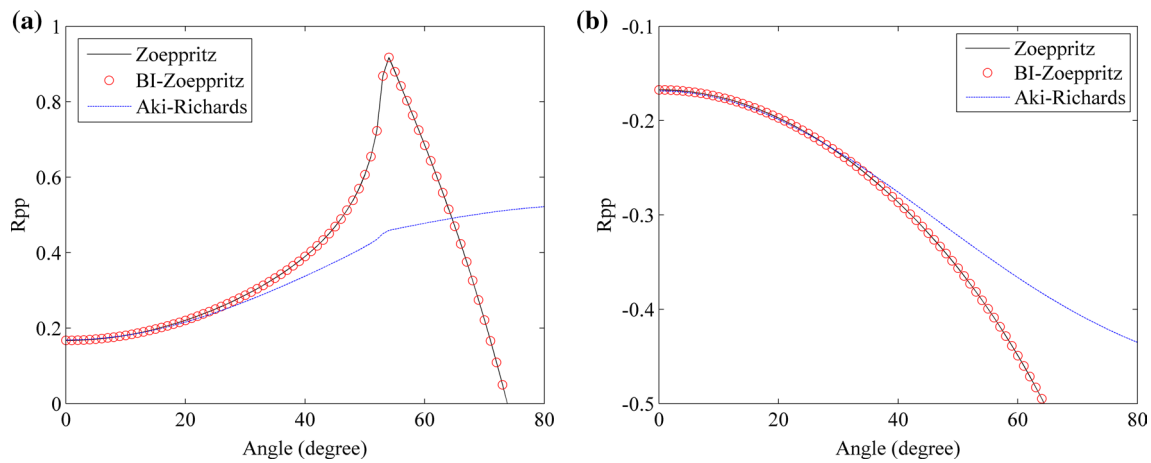
### AVO inversion for BI

Prestack AVO inversion is a highly nonlinear optimization problem, which can be generally solved by weighted

**Table 1** Elastic parameter of Ostrander shale model

Lithology	$V_p$ (m/s)	$V_s$ (m/s)	$\rho$ (g/cc)	YM ( $\text{kg m}^{-1} \text{ s}^{-2}$ )	PR	BI
Sandstone	2438	1625	2.14	$1.0400 \times 10^{10}$	0.4001	68.7656
Shale	3048	1244	2.40	$1.2435 \times 10^{10}$	0.1003	15.1031
Sandstone	2438	1625	2.14	$1.0400 \times 10^{10}$	0.4001	68.7656





**Fig. 1** A comparison of  $R_{pp}$  between the exact Zoeppritz equation, BI\_Zoeppritz equation and Aki–Richards approximation. **a** Interface between sandstone and shale; **b** interface between shale and sandstone

coefficient method for objective function. However, this method requires artificially setting a set of weight coefficients to form a single objective function. Therefore, the influence of weight coefficients on the inversion results is inevitable. Another effective method to solve the inversion problem is the multi-objective function optimization algorithm, which can maximize or minimize the multi-objective function under some constraints in the decision space (Liu and Wang 2018). It avoids the subjective influence of choosing weight coefficient by human. The main purpose of multi-objective function optimization algorithm is to obtain a set of Pareto optimal solutions which satisfies the multi-objective function at the same time. The image of Pareto optimal solution set in objective space is called Pareto optimal frontier. NSGA II is a global optimization algorithm with elite strategy, which can minimize multi-objective function without setting the weight coefficient artificially (Deb et al. 2002), so as to extract the Pareto optimal solution according to the needs of the actual problem.

To obtain the  $V_p$ ,  $V_s$  and BI from seismic records by NSGA II, the cross-correlation principle (Li and Mallik 2015) and the least root mean square error principle are employed to construct the multi-objective function as follows:

$$\begin{cases} \text{CCE} = 1 - \frac{S_{pp} \cdot D_{pp}}{\sqrt{S_{pp} \cdot S_{pp}} \cdot \sqrt{D_{pp} \cdot D_{pp}}} \\ \text{RMSE} = \sqrt{\frac{1}{m} \sum_{i=1}^m (S_{pp} - D_{pp})^2}, \quad i = 1, 2, \dots, m \end{cases}, \quad (13)$$

where  $D_{pp}$  is the actual observed PP-wave seismic record,  $S_{pp}$  is the PP-wave synthetic record obtained from forward modeling of elastic parameters by using the BI\_Zoeppritz equation,  $m$  is the sampling numbers, RMSE represents

root mean square error reflecting the difference between the observed seismic record and synthetic record, and CCE represents cross-correlation error reflecting the correlation between them, respectively. We use CCE and RMSE as a multi-objective function to obtain more accurate inversion results.

In this study, we employ NSGA II to get the solution of the multi-objective function shown in Eq. 13 and perform the AVO inversion for BI. We consider the potential solution of the elastic parameters ( $V_p$ ,  $V_s$  and BI) as an individual in the  $D$ -dimensional solution space. For the  $i$ th individual, NSGA II maps a set of solutions  $x^i = [x_1^i, x_2^i, \dots, x_D^i]^T$  in solution space to the multi-objective function values  $y^i = [y_1^i, y_2^i]^T$  in objective space to find a set of Pareto optimal solutions that can satisfy the multi-objective function at the same time. The overall process is as follows:

- (1) Generate a population of  $N$  individuals and obtain the multi-objective function value of each individual by using Eq. 13;
- (2) Use a linear transformation to get the corresponding transformed multi-objective function value of each individual;
- (3) Assign two attribute parameters of nondominated level and crowding distance to each individual;
- (4) Apply a series of operations, such as tournament selection, SBC and RPM, to get the offspring population of size  $N$ ;
- (5) Merge the parent and the offspring populations into a single population with  $2N$  individuals;
- (6) Select  $N$  new individuals from the single population by nondominated sorting and crowding distance calculation;
- (7) Arrange the new individuals as parent population to participate in the next genetic operation;

- (8) Repeat (4) to (7) until the preset stopping criteria are met;
- (9) Obtain the Pareto optimal solution set of the last generation population;
- (10) Extract the optimal elastic parameters ( $V_p$ ,  $V_s$  and BI) from the Pareto optimal solution set.

NSGA II algorithm mainly includes two mechanisms: fast nondominated sorting and population diversity protection. Assuming that there is a population with  $N$  individuals, the objective function value of each individual can be calculated by Eq. 13. Then, each individual is assigned a nondominated level according to the objective function value, which is called nondominated sorting. Nondominated sorting follows the rule: individuals of nondominated level 1 dominate individuals of other nondominated levels, but the internal individuals of same nondominated level do not dominate each other; Individuals of nondominated level 2 dominate individuals of nondominated level 3 and above. Similarly, there is no domination between the individuals of same nondominated level, and so on until all individuals of the population are assigned to a nondominated level. At the same time, in order to prevent the population from losing diversity in the process of genetic operation, NSGA II assigns the crowding distance parameter to the individuals of same nondominated level. To calculate the crowding distance, individuals of the same nondominated level are sorted in ascending order by objective function value, and then the individuals at both ends of the order are assigned an infinite number as their crowding distance, while other individuals can calculate different crowding distances according to certain criteria. The individuals with larger crowding distance in the same nondominated level are more likely to be selected to participate in the next genetic operation. The main purpose of these two mechanisms is to search for the solution set that is close to the real Pareto optimal solution set, while maintaining population diversity. Deb et al. (2002) described the detailed steps of traditional NSGA II algorithm. To ensure the accuracy of AVO inversion, Liu and Wang (2018) improved the implement of NSGA II and obtained good inversion results.

The traditional NSGA II algorithm performs nondominated sorting and crowding distance calculation based on the original objective function value of each individual. However, since most individuals in the early population of genetic iterations are far away from the global optimal solution, the two mechanisms based on the original objective function value often cause the superior individuals to drive the traditional NSGA II algorithm to fall into the local minimum (Goldberg 1989). To avoid this phenomenon, the original objective function value is processed by linear transformation to reduce the fitness value of the superior individual and increase the fitness value of the inferior individual, so

that the inversion process of multi-objective function stably converges to the global optimal direction.

The linear transformation is given in Eqs. 14 and 15:

$$f' = af + b, \quad (14)$$

and

$$\begin{cases} a = \frac{f_{\text{avg}}(C_m - 1)}{f_{\text{max}} - f_{\text{avg}}} \\ b = f_{\text{avg}}(1 - a) \\ f'_{\text{max}} = C_m f_{\text{max}} \\ C_m = 1.2 + 0.8 \times \frac{t}{t_{\text{max}}} \end{cases}, \quad (15)$$

where  $f$  and  $f'$  represent the value of original objective function and corresponding transformed objective function,  $f_{\text{avg}}$  and  $f_{\text{max}}$  represent the average and maximum values of original objective function,  $t$  and  $t_{\text{max}}$  represent the current and maximum value of genetic iterations, respectively.  $C_m$  is a parameter changing gradually with genetic iterations, and  $f'_{\text{max}}$  represent the maximum value of the transformed objective function.

In order to find the global optimal solution of multi-objective function, simulated binary crossover (SBC) and real parameter mutation (RPM) operators are used to continuously sample in decision space by real value coding.

In SBC, first, a parameter  $\beta$  is defined as (Deb and Agrawal 1995):

$$\beta = 1 + \frac{2}{x_i^{2,t} - x_i^{1,t}} \min[(x_i^{1,t} - x_i^{L,t}), (x_i^{U,t} - x_i^{2,t})], \quad (16)$$

where  $x_i^{1,t}$  and  $x_i^{2,t}$  represent the  $i$ -th coding value of the parent individual 1 and 2 in the  $t$ th generation,  $x_i^{L,t}$  and  $x_i^{U,t}$  represent the minimum and maximum coding values for individuals in the  $t$ th parent population, respectively. Then, a parameter  $\alpha$  is defined as:

$$\alpha = 2 - \frac{1}{\beta^{\eta+1}}, \quad (17)$$

where a nonnegative real number  $\eta$  denotes the cross-distribution index. The larger  $\eta$  value is, the closer the offspring is to its parent. A parameter  $\beta_{qi}$  is given as:

$$\beta_{qi} = \begin{cases} (u_i \alpha)^{\frac{1}{\eta+1}}; & u_i \leq \frac{1}{\alpha} \\ \left(\frac{1}{2-u_i \alpha}\right)^{\frac{1}{\eta+1}}; & \text{otherwise} \end{cases}, \quad (18)$$

where  $u_i$  represents a random number between 0 and 1. Finally, we can use Eq. 19 to obtain two offspring individuals:

$$\begin{cases} x_i^{1,t+1} = 0.5[(x_i^{1,t} + x_i^{2,t}) - \beta_{qi}(x_i^{2,t} - x_i^{1,t})] \\ x_i^{2,t+1} = 0.5[(x_i^{1,t} + x_i^{2,t}) + \beta_{qi}(x_i^{2,t} - x_i^{1,t})] \end{cases}, \quad (19)$$

where  $x_i^{1,t+1}$  and  $x_i^{2,t+1}$  represent the  $i$ th coding value of the offspring individuals 1 and 2 in any generation  $t$ , respectively.

In order to avoid the genetic operation falling into local minimum, we do real parameter mutation operation on the offspring individuals generated by SBC.

In RPM, first, a parameter  $\delta$  is defined as (Deb and Agrawal 1999):

$$\delta = \frac{\min[(x_i^{\text{child}} - x_i^L), (x_i^U - x_i^{\text{child}})]}{x_i^U - x_i^L}, \quad (20)$$

where  $x_i^{\text{child}}$  represents the  $i$ th coding value of the offspring individual after SBC,  $x_i^L$  and  $x_i^U$  represent the  $i$ th minimum and maximum coding values of the offspring individuals after SBC, respectively. Parameter  $\bar{\delta}_i$  can be expressed as:

$$\bar{\delta}_i = \begin{cases} [2r_i + (1 - 2r_i)(1 - \delta)^{k+1}]^{\frac{1}{k+1}} - 1; & r_i \leq 0.5 \\ 1 - [2(1 - r_i) + 2(r_i - \frac{1}{2})(1 - \delta)^{k+1}]^{\frac{1}{k+1}}; & \text{otherwise} \end{cases}, \quad (21)$$

where  $r_i$  represent a random number between 0 and 1, a non-negative real number  $k$  denotes the mutation distribution index. The larger  $k$  value is, the closer the offspring is to its parent. Then, we can obtain the offspring individuals after mutation according to Eq. 22:

$$x_i^{\text{mutated}} = x_i^{\text{child}} + (x_i^U - x_i^L)\bar{\delta}_i, \quad (22)$$

where  $x_i^{\text{mutated}}$  is the  $i$ th coding value of the offspring individual after mutation.

In order to stabilize the convergence process of AVO inversion, we set four key parameters of SBC and RPM to change with the genetic iterations as (Li and Mallick 2015; Liu and Wang 2018):

$$\begin{cases} P_c = 0.7 - 0.1 \times \frac{t}{t_{\text{max}}} \\ P_m = \frac{1}{n} + \frac{t}{t_{\text{max}}} \times \left(1 - \frac{1}{n}\right), \\ \eta = 1.0 + 19.0 \times \frac{t}{t_{\text{max}}} \\ k = 100 + t \end{cases}, \quad (23)$$

where  $n$  represents the total number of variables,  $P_c$  and  $\eta$  represent crossover probability and crossover distribution index,  $P_m$  and  $k$  represent mutation probability and mutation distribution index, respectively.

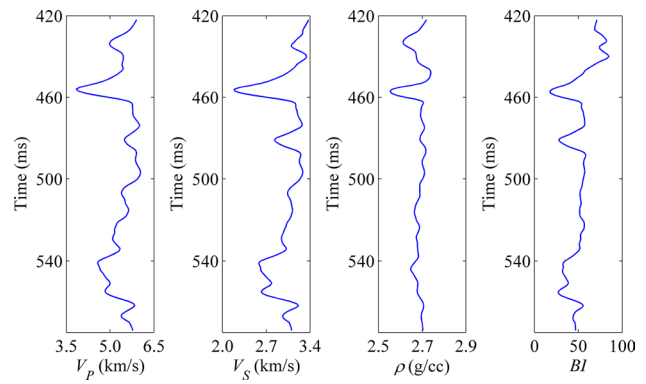


Fig. 2 Elastic parameters of theoretical model

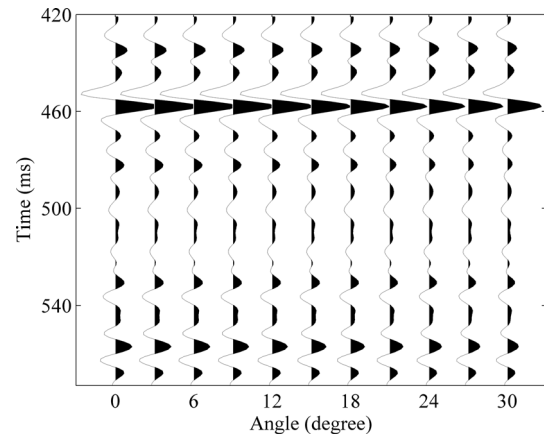


Fig. 3 PP-wave angle gathers of theoretical model without noise

## Model data test

### Model test I

To verify the validity of the inversion method, we select the actual logging data of a shale gas field in China, and build a one-dimensional theoretical model as shown in Fig. 2 by smooth filtering. The BI of theoretical model is obtained by Eq. 6, and the  $R_{pp}$  is obtained by the BI\_Zoeppritz equation. Based on the convolution of  $R_{pp}$  and a 35-Hz Ricker wavelet, the synthetic angle gathers for an incident angle of  $0^\circ$ – $30^\circ$  are generated as shown in Fig. 3.

We apply the AVO inversion to the theoretical model and design an initial model (green line shown in Fig. 4) through low-pass filtering to improve the inversion accuracy. Two different search windows (black line shown in Figs. 5 and 6) are also used to constrain the inversion process and test the impact of the search window on inversion accuracy. One is a linear window with constant maximum and minimum.

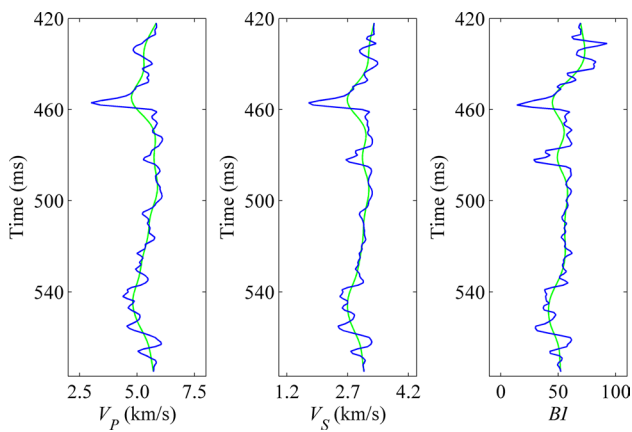


Fig. 4 Initial model for theoretical model

And another is a optimized window referred to the initial model. In addition, the population size and total number of genetic iterations of the NSGA II are set to 800 and 2000, respectively. Figure 5a and b shows the inverted elastic parameters using the optimized window and linear window, respectively. Figure 6 displays the corresponding inversion errors of two different search windows. From Figs. 5 and 6, we can find that the inverted elastic parameters are in good agreement with the real values of theoretical model. However, some large errors occur above 540 ms, especially at the positions with violent variation of elastic parameters. Additionally, the quality of search window has a big impact on inversion accuracy. Compared with the linear search window, we can obtain better inversion results by using the optimized search window.

To test the noise immunity of the inversion method, we add random noise with a signal-to-noise ratio (SNR) of 3 on

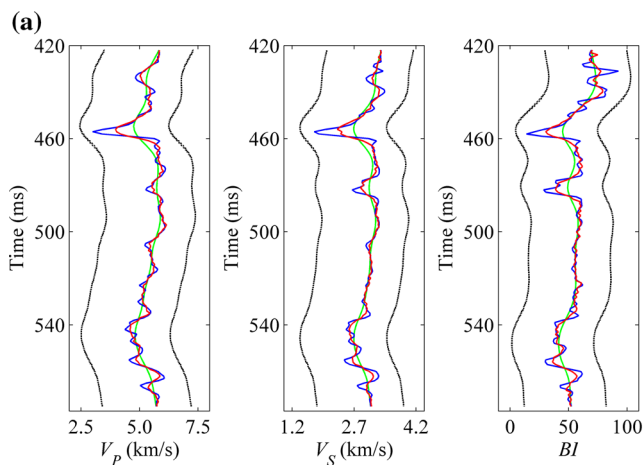


Fig. 5 Inverted elastic parameters without noise using a the optimized search window and b the linear search window. The green line and blue line denote the initial model and theoretical model, and the

the angle gathers, as shown in Fig. 7. The linear search window and optimized search window are employed to constrain the inversion process. Figure 8a and b displays the inversion results using the optimized window and linear window, respectively. Figure 9 shows the corresponding absolute errors of inversion results. The annotation in Figs. 8 and 9 is the same as Figs. 5 and 6, respectively. From these figures, we can see that the random noise has a big impact on the inversion accuracy, especially for the BI using linear search window. When the noise is added, the inversion accuracy of linear search window reduces obviously, whereas the inversion accuracy of optimized search window changes little.

Figure 10 shows the multi-objective function values (CCE and RMSE) of different generations in four cases. Figure 11 gives the population distributions of the last generation in four cases, which reflect the Pareto optimal front. Note from

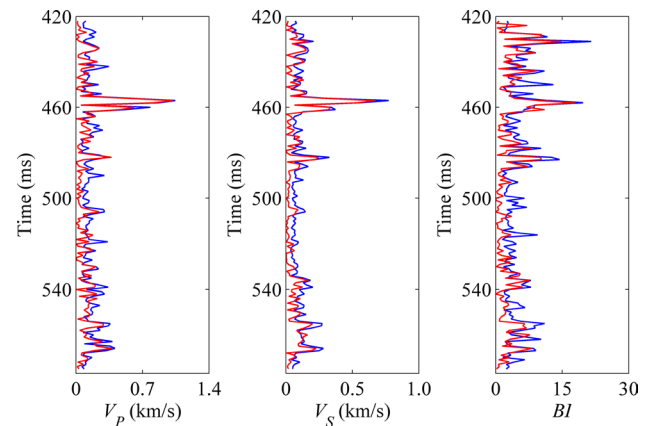
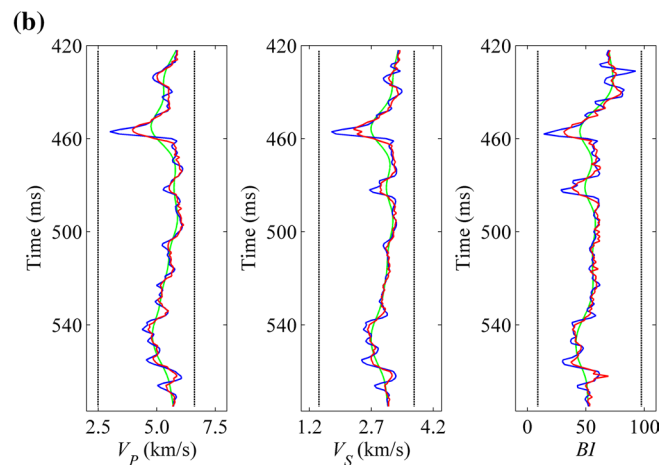


Fig. 6 Error comparison of inverted elastic parameters without noise using the optimized search window (red line) and the linear search window (blue line)



black line and red line represent the search windows and corresponding inverted elastic parameters

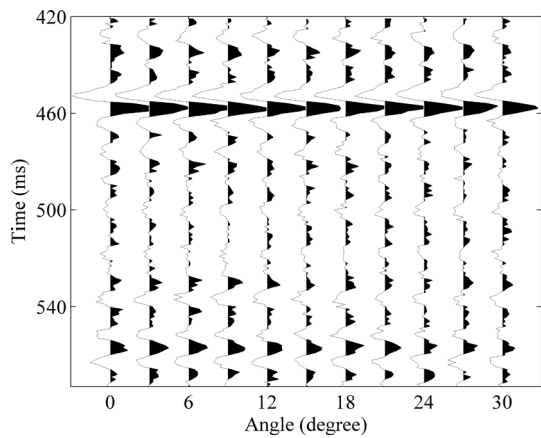


Fig. 7 PP-wave angle gathers of theoretical model with SNR = 3

Figs. 10 and 11, random noise and search window have certain influence on the convergence of the AVO inversion. The convergence degree of CCE, RMSE and Pareto optimal front will be worse when the search window moves away from the trend of theoretical model. Similarly, when the angle gathers are added with random noise, the convergence degree of CCE, RMSE and Pareto optimal front also deteriorate, which is more obvious on RMSE. In the inversion process of one-dimensional theoretical model, the inversion method converges quickly in the first 200 genetic iterations, then slows down with the increase in genetic iterations and keeps stable at last. Therefore, the appropriate number of genetic iterations can significantly reduce the computation cost while ensuring the convergence degree and inversion accuracy.

### Model test II

In order to make the test of theoretical model shown in Fig. 2 more close to the actual circumstances, we use a 30-Hz Ricker wavelet to generate the synthetic angle gathers for an incident angle of 0°–30°. Some random noise with SNR = 1 is added to the angle gathers. The optimized search window is also employed to constrain the inversion process. Figure 12a and b displays the inversion results without noise and with SNR = 1, respectively. Figure 13 shows the corresponding absolute errors of inversion results. From these figures, we can see that the inversion results without noise are in good agreement with the real values of theoretical model. Compared with the inversion results using a 35-Hz Ricker wavelet, the absolute errors of the inversion result using a 30-Hz Ricker wavelet are not very big. When the

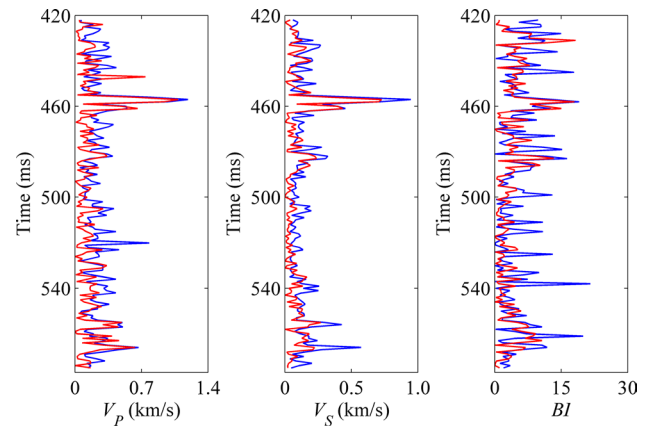


Fig. 9 Error comparison of inverted elastic parameters with SNR = 3 using the optimized search window (red line) and the linear search window (blue line)

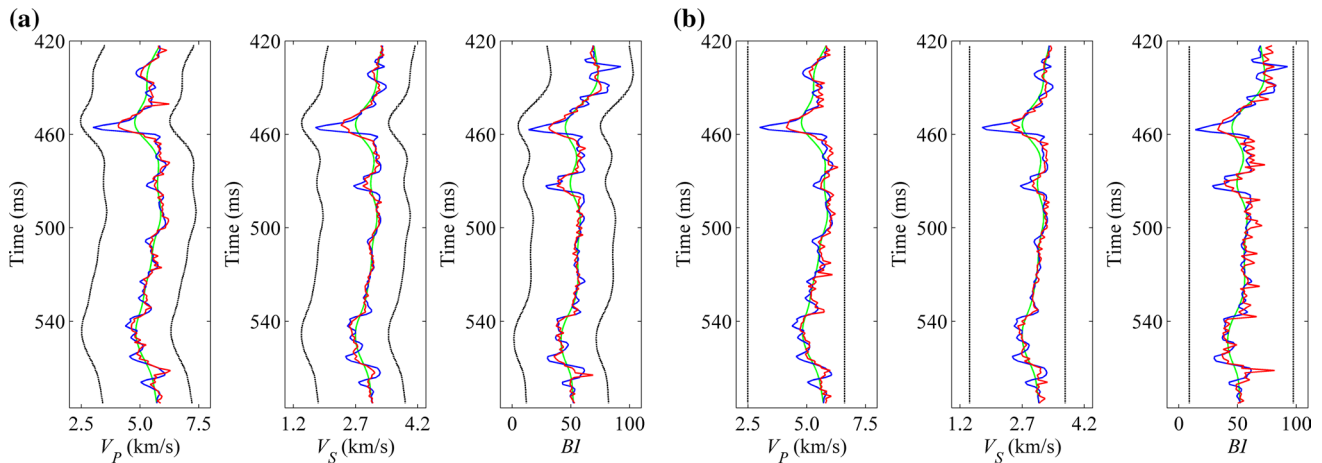
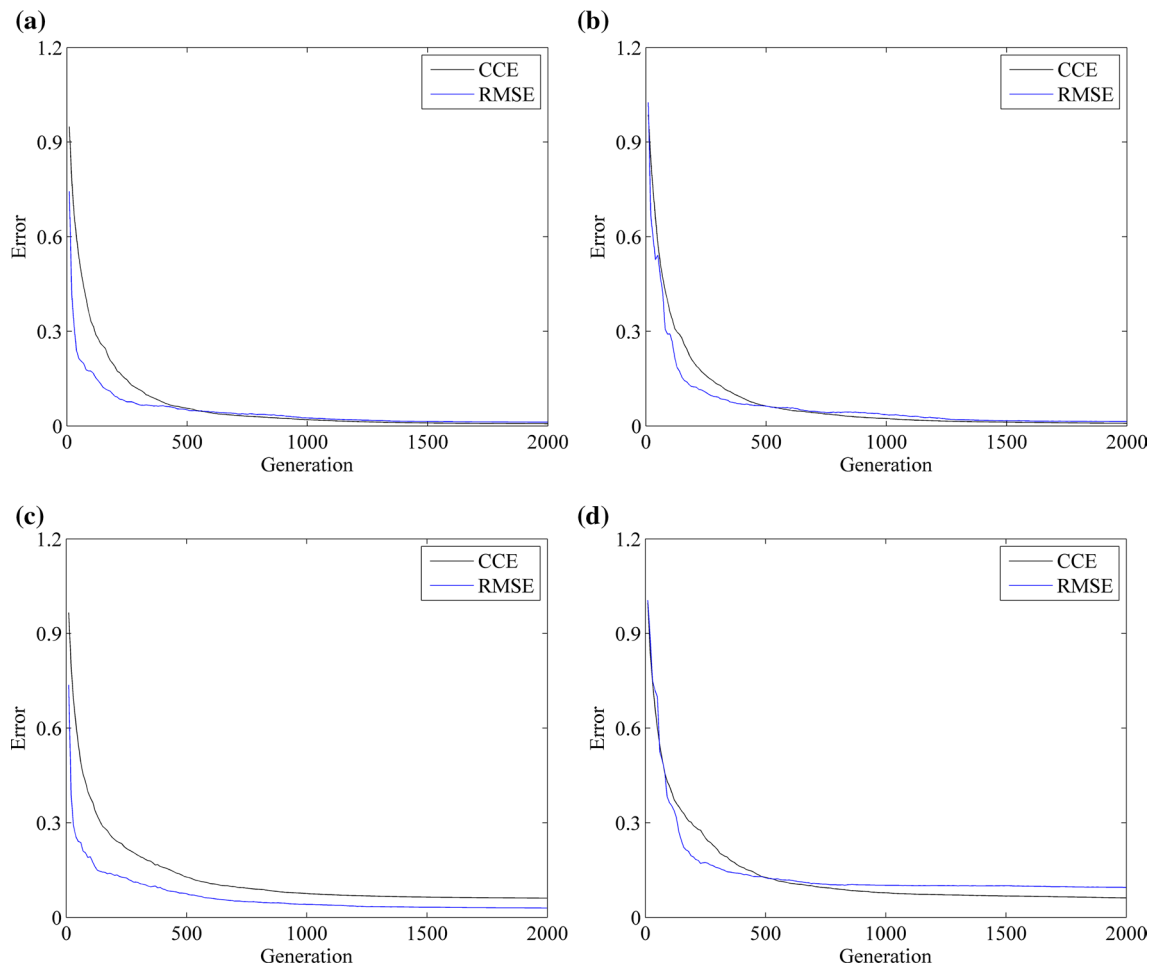


Fig. 8 Inversion results with SNR = 3 using **a** the optimized search window and **b** the linear search window. The green line and blue line denote the initial model and theoretical model, and the black line and red line denote the search windows and corresponding inverted elastic parameters



**Fig. 10** Comparison of convergence speed in four cases: **a** optimized search window and **b** linear search window for angle gathers without noise, and **c** optimized search window and **d** linear search window for angle gathers with SNR=3

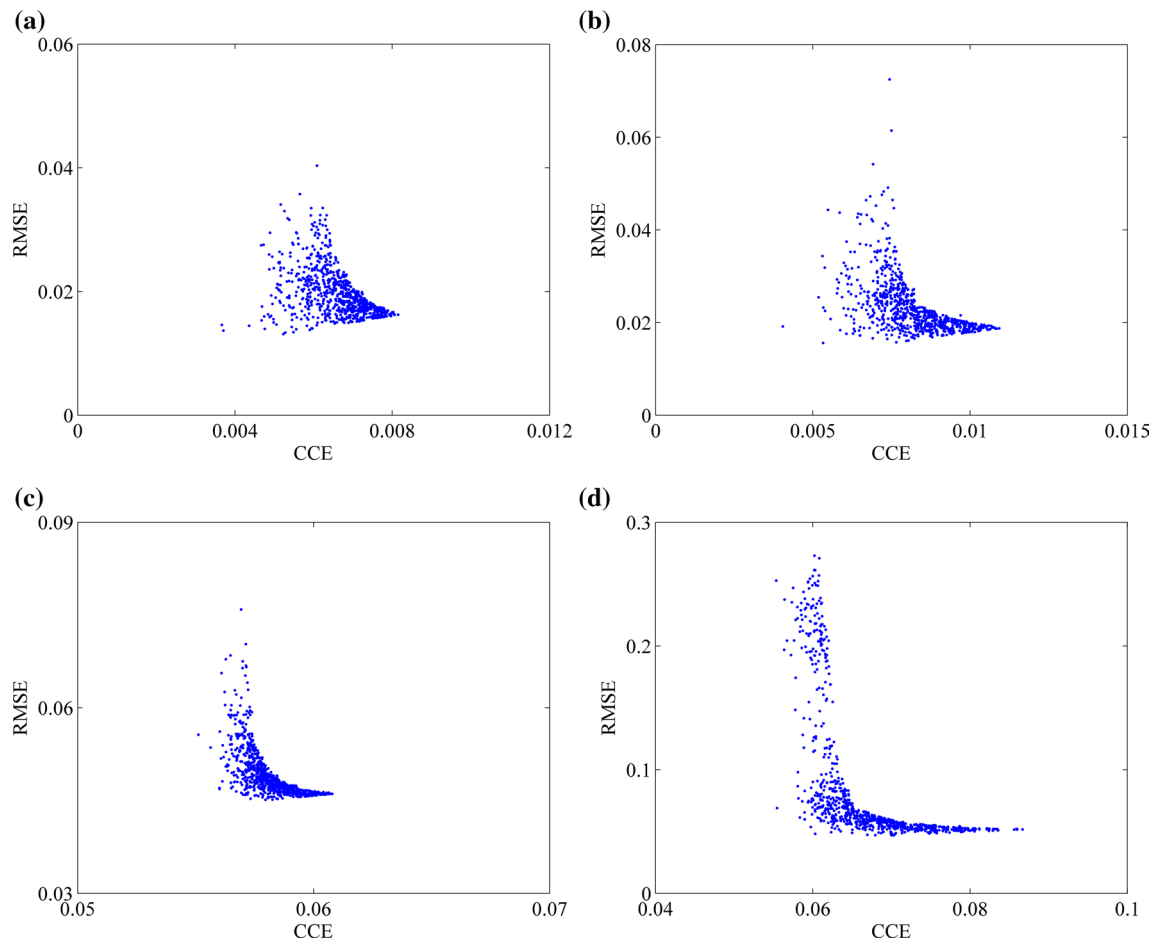
noise with SNR = 1 is added, the inversion accuracy reduces obviously and the absolute errors of three elastic parameters show a random distribution. From Figs. 9 and 13, we can see that the random noise has a big impact on the inversion accuracy. Nevertheless, the inverted elastic parameters still accord with the theoretical model, especially under the condition of low SNR. Overall, the proposed inversion method using the optimized search window has high precision and high noise immunity.

## Actual data application

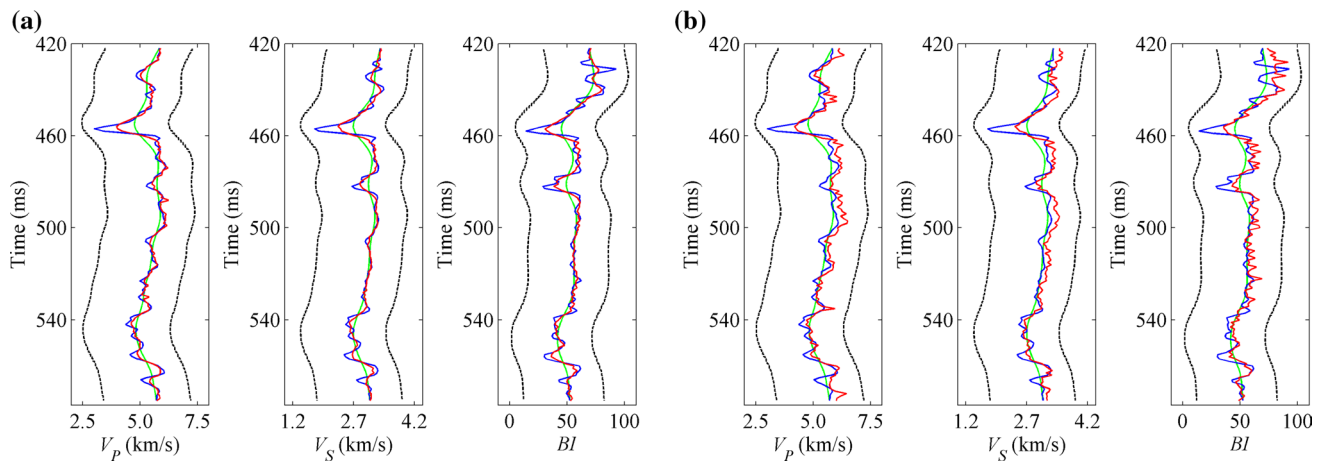
### One-dimensional data application

To further verify the applicability and effectiveness of the inversion method, we select actual logging data from a shale gas field in southern China to test it. Figure 14 shows the depth-domain logging curves of actual data. By using the

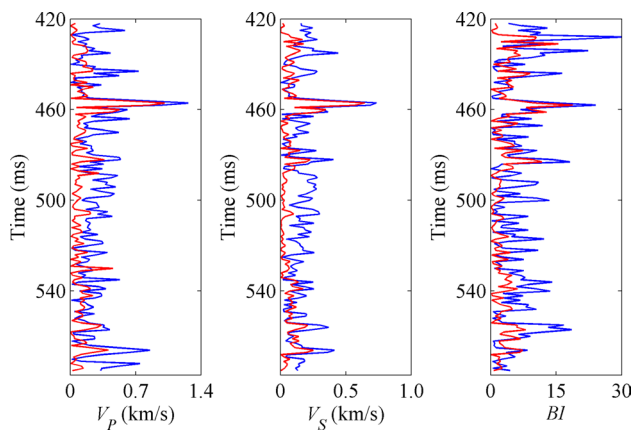
kelly bushing, seismic datum and replacement velocity, the depth-domain logging curves are converted to the time domain. Figure 15 gives the time domain logging curves (blue line) with a sampling interval of 1 ms, and there are 369 sampling points in total. Based on a 35-Hz Ricker wavelet and BI\_Zoeppritz equation, the PP-wave synthetic angle gathers with an incident angle of  $0^{\circ}$ – $30^{\circ}$  are generated as shown in Fig. 16a. Random noise with SNR of 3 is added on the angle gathers, as shown in Fig. 16b. The goal of the inversion method is to obtain  $V_p$ ,  $V_s$  and BI, so the corresponding model parameter vector has 1107 unknowns. To improve the inversion accuracy, an initial model, an optimized search window and an optimized initial population are introduced into the inversion process. The green line in Fig. 15 denotes the initial model, which is the result of smoothing the actual logging data. The optimized search window is the search range obtained by expanding the initial model to both sides and including all possible values of the elastic parameters. The optimized initial population is the



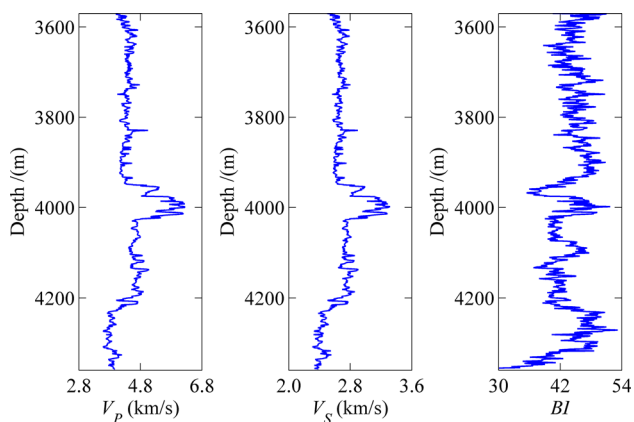
**Fig. 11** Population distributions at the last generation in four cases: **a** optimized search window and **b** linear search window for angle gathers without noise, and **c** optimized search window and **d** linear search window for angle gathers with SNR = 3



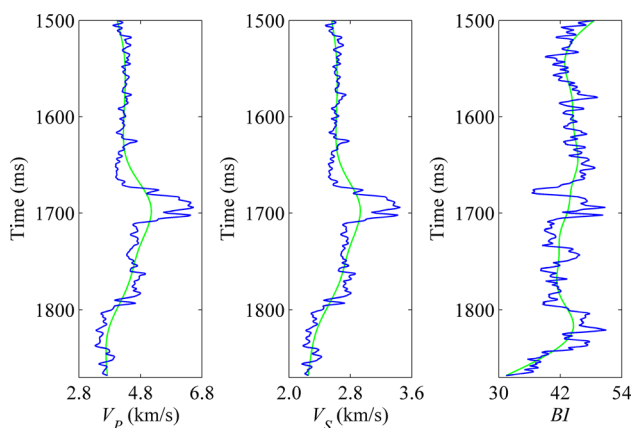
**Fig. 12** Inversion results of model data **a** without noise and **b** with SNR = 1. The green line and blue line denote the initial model and theoretical data, and the black line and red line denote the search windows and corresponding inverted elastic parameters



**Fig. 13** Error comparison of inverted elastic parameters without noise (red line) and with SNR = 1 (blue line)



**Fig. 14** Actual data in depth domain



**Fig. 15** Actual data in time domain. The blue line and green line represent the logging data and corresponding initial model, respectively

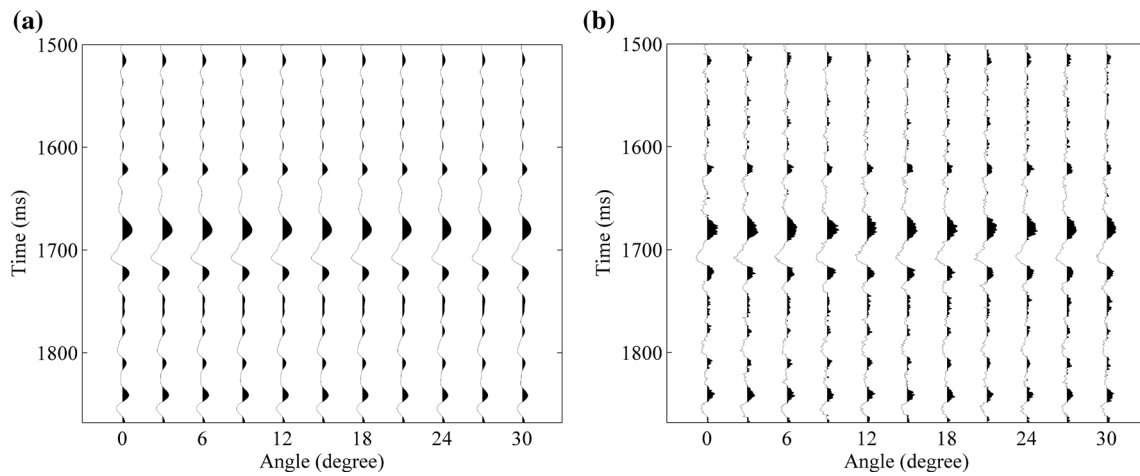
population of  $N$  mutant individuals in the initial model. In the inversion process, we set the population size as 800 and the maximum number of genetic iterations as 2000.

Figure 17a gives the inverted elastic parameters without noise. It can be noted from Fig. 17a that the inversion method can get more reasonable  $V_p$  and  $V_s$ . However, the inversion error of BI is small only in the regions where the BI changes slowly. This is because the far-angle gathers are not used in the process of AVO inversion. Compared with low-angle gathers, wide-angle gathers can be used to obtain more accurate elastic parameters (Virieux and Operto 2009). Accordingly, it is not surprising that the inversion results of BI around 1680 ms and 1835 ms are poor. Figure 17b gives the inverted elastic parameters with SNR = 3. Figure 18 gives the corresponding absolute error. The lines in red denote the error between the inverted elastic parameters without noise and the actual data, and the lines in blue indicate the error between the inverted elastic parameters with SNR = 3 and the actual data. From Fig. 18, we can see that the inversion accuracy decreases when the angle gathers contain random noise (SNR = 3). This is more evidence for the fact that the noise has a great influence on inversion accuracy. Nevertheless, the inverted elastic parameters still accord with the actual data.

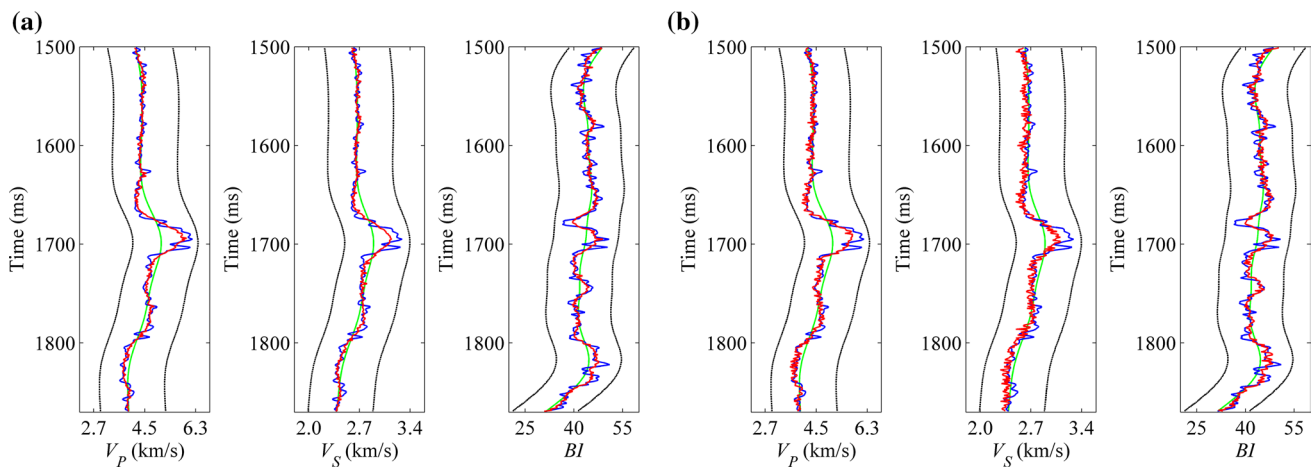
Figures 19a and 20a show the synthetic angle gathers computed by using the inversion results of Fig. 17a and b, respectively. Figure 19b gives the difference section between the synthetic angle gathers of Fig. 19a and the raw angle gathers without noise of Fig. 16a. Figure 20b gives the difference section between the synthetic angle gathers of Fig. 20a and the raw angle gathers with SNR = 3 of Fig. 16b. From Figs. 19 and 20, we can find that the synthetic angle gathers calculated by the inversion results are close to the raw angle gathers without noise as shown in Fig. 16a. Only a few weak events and random noises remain on the different sections, which proves that the inversion method has good feasibility and noise immunity.

Figure 21 displays the evolution of the initial population with genetic iterations. The horizontal and vertical axes represent CCE and RMSE, respectively, which are the multi-objective function values after linear transformation. From this figure, we can see that all individuals of the initial population are arbitrarily and widely distributed in the objective space and then gradually converge to the global optimal solution, which denotes that NSGA II can effectively deal with the multi-objective function optimization problem. In addition, the inversion method converges quickly in the early stage of genetic iterations and then slows down with the increase in genetic iterations. Figure 22 gives the

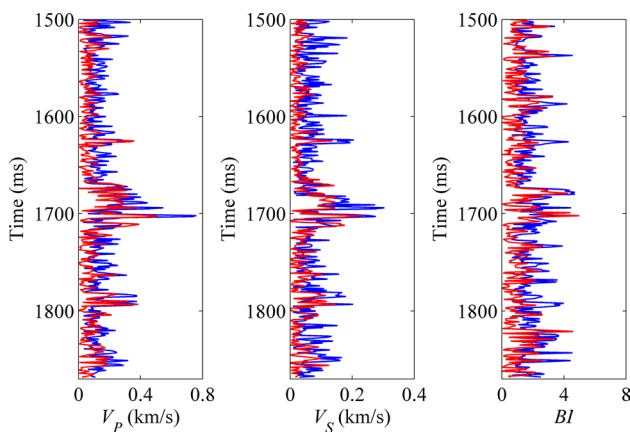




**Fig. 16** PP-wave angle gathers **a** without noise and **b** with SNR = 3



**Fig. 17** Inversion results of actual data **a** without noise and **b** with SNR = 3. The green line and blue line denote the initial model and actual data, and the black line and red line denote the search windows and corresponding inverted elastic parameters

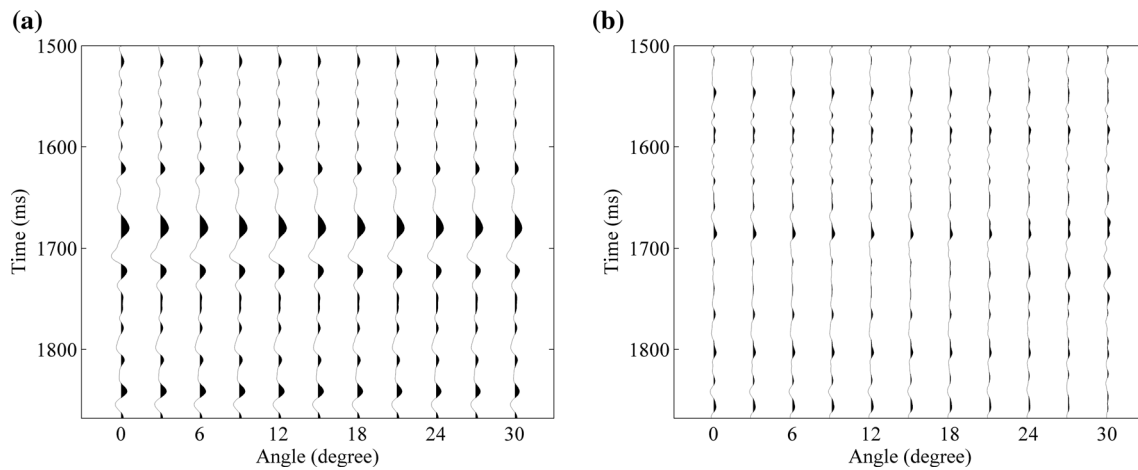


**Fig. 18** Error comparison of inverted elastic parameters without noise (red line) and with SNR = 3 (blue line)

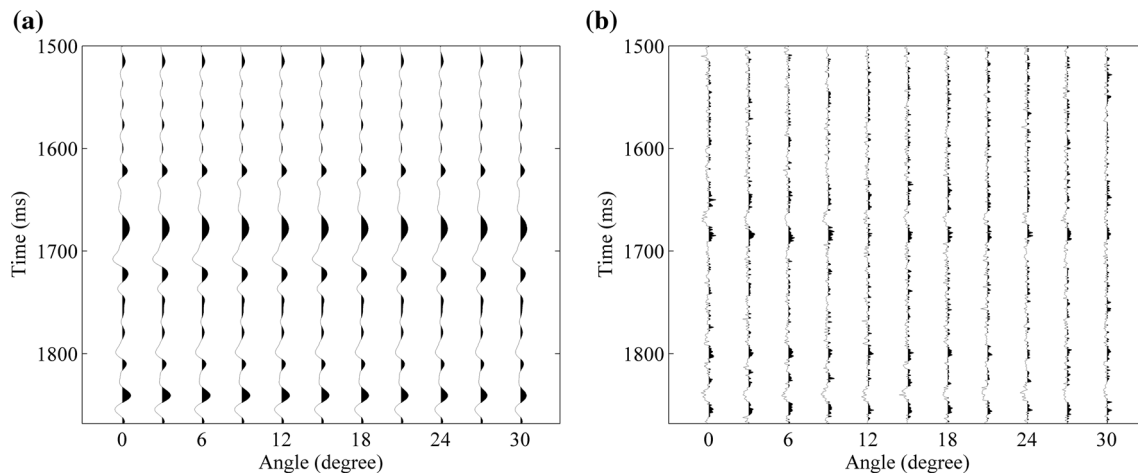
population distribution of the last generation in the objective space. It can be seen from Fig. 22 that Pareto optimal front also has poor convergence because of the random noise in angle gathers.

**Two-dimensional data application**

We use actual seismic data obtained from a shale gas fields in southwest China to test the proposed method. To suppress noise and improve the SNR of seismic data, partial stack of incident angle gathers is usually needed before inversion. Figure 23 shows a two-dimensional (2D) section of the PP-wave seismic data based on partial angle stack ranges of 4°–16°, 15°–25° and 24°–36°. The main frequencies of three partial angle stack data are 31 Hz, 30 Hz and 29 Hz,



**Fig. 19** **a** Synthetic angle gathers for inverted elastic parameters of Fig. 17a, **b** the difference section between the synthetic angle gathers and raw angle gathers without noise of Fig. 16a

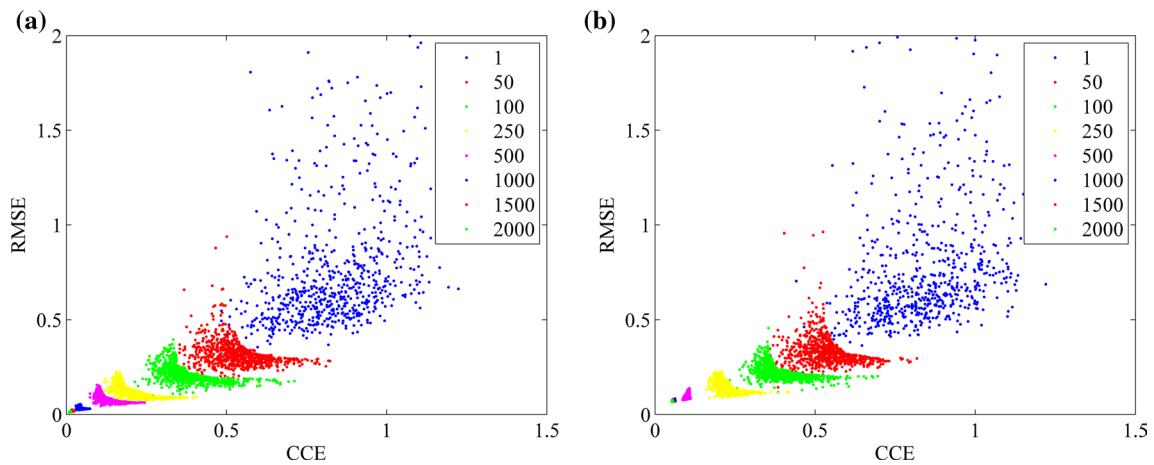


**Fig. 20** **a** Synthetic angle gathers for inverted elastic parameters of Fig. 17b, **b** the difference section between the synthetic angle gathers and raw angle gathers with SNR=3 without noise of Fig. 16b

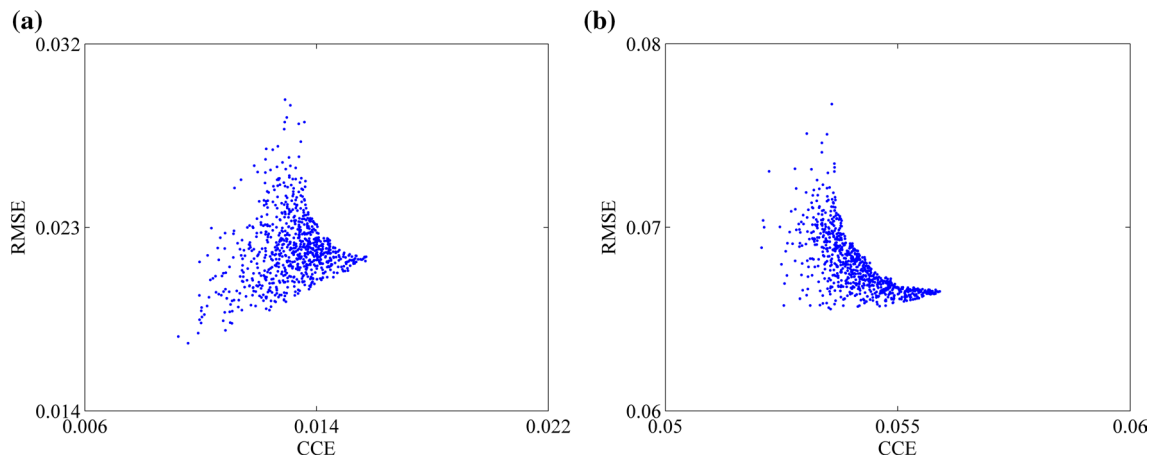
respectively. Well C is located in CDP 1528. The YM and PR of Well C are calculated by using Eqs. 4 and 5. The statistical analysis indicates that  $YM_{\min} = 1.2365 \times 10^{10} \text{ kg m}^{-1} \text{ s}^{-2}$ ,  $YM_{\max} = 8.0345 \times 10^{10} \text{ kg m}^{-1} \text{ s}^{-2}$ ,  $PR_{\min} = 0.1409$  and  $PR_{\max} = 0.3667$ . Therefore, the corresponding BI can be calculated according to Eq. 6. To perform well-seismic calibration, a 30-Hz Ricker wavelets is created and the  $R_{PP}$  is calculated by using BI\_Zoeppritz equation. To be consistent with the actual seismic data, the PP-wave synthetic angle gathers are generated and partially stacked according to angle ranges of  $4^{\circ}$ – $16^{\circ}$ ,  $15^{\circ}$ – $25^{\circ}$  and  $24^{\circ}$ – $36^{\circ}$ . Based on the well-seismic calibration, the time domain well logs of Well C are shown as the blue line in Fig. 25.

In order to perform the AVO inversion for BI of the 2D actual seismic data, we establish a 2D initial model and a

2D optimized search window based on horizons and logging data of Well C. In Fig. 25, the green line and black line denote the initial model and the optimized search windows for the seismic trace near the borehole (CDP 1528). In the inversion process, we set the population size as 800 and the maximum number of genetic iterations as 2000 for each trace of the 2D seismic section. Figure 24 shows the 2D elastic parameters section inverted by the proposed method. To further observe the inversion accuracy of the proposed method, the red line in Fig. 25 gives the inverted results of the seismic trace near the borehole (CDP 1528). It can be seen that the inversion results of the proposed method are in good agreement with actual logging data, especially near the 1680–1740 ms. In summary, it is feasible to obtain BI from actual seismic data by using the proposed method.



**Fig. 21** Evolution of the initial population with genetic iterations using angle gathers **a** without noise and **b** with SNR=3. The different colored dots represent individuals of different genetic iterations



**Fig. 22** Cross-plot of CCE and RMSE at the last generation using angle gathers **a** without noise and **b** with SNR=3

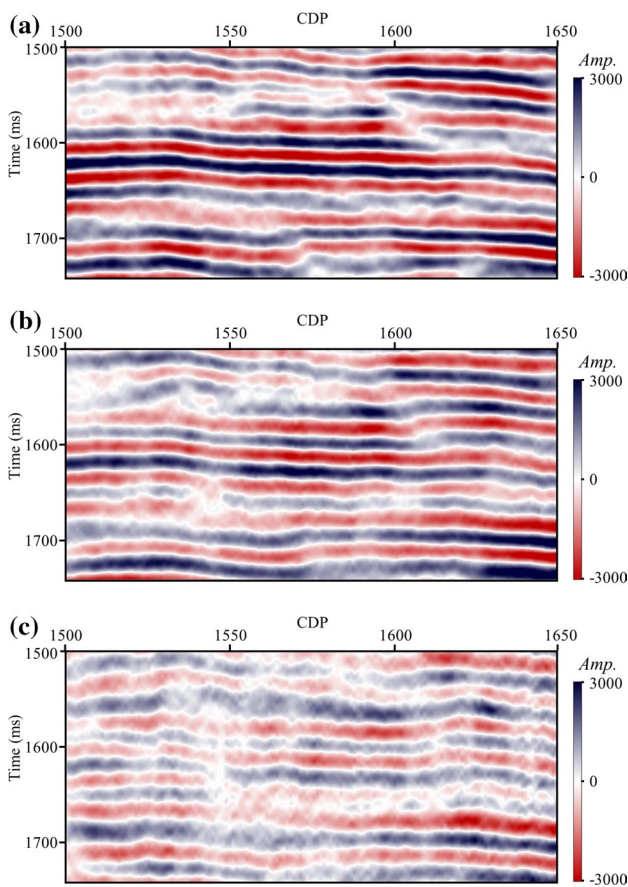
## Discussion

Because NSGA II is a global optimization algorithm for multi-objective function, the proposed inversion method in this paper can estimate the global optimal solution of  $V_p$ ,  $V_s$  and BI. However, the computation cost of NSGA II is very large due to the iterative calculation of multi-objective function values by using the BI-Zoeppritz equation. Therefore, the maximum value of genetic iterations cannot be set too large.

For complex geological models with drastic changes in elastic parameters, the proposed inversion method can constrain the entire inversion process with an initial model and an optimized search window to reduce the calculation time and improve the inversion accuracy. In the inversion of actual seismic data, the initial model of  $V_p$  and  $V_s$  can also be constructed by velocity analysis in seismic data processing. The initial model of BI can also be defined by

rock physics. An initial population consisting of  $N$  mutant individuals can be generated to approximate the real model through nondominated sorting, tournament selection, simulated binary crossover, real parameter mutation and elitism (Liu and Wang 2018).

Although the BI with high accuracy can be obtained by the proposed inversion method, there are two factors that affect the inversion accuracy. One is the range of elastic parameters. Figure 4 shows the  $V_p$  (2.5–7.5 km/s),  $V_s$  (1.2–4.2 km/s), and BI (0–100) of theoretical model, while Fig. 15 shows the  $V_p$  (2.8–6.8 km/s),  $V_s$  (2.0–3.6 km/s) and BI (30–54) of actual data. By comparing Fig. 6 with Fig. 18, we can find that the larger the variation range of elastic parameters, the larger the error of inversion results. The other is noise that increases the uncertainty of inversion results, so noise should be suppressed as much as possible in seismic data processing.

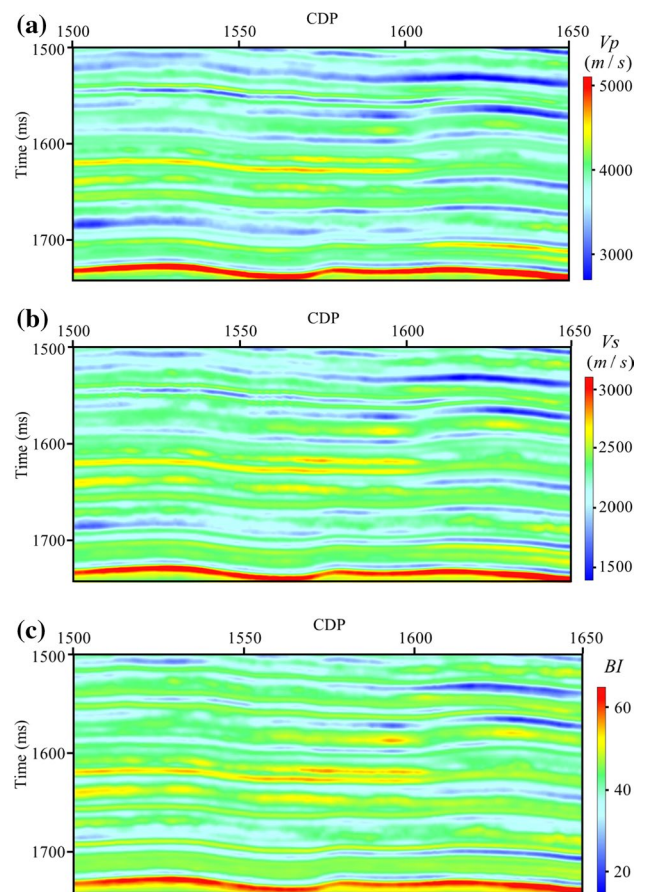


**Fig. 23** PP-wave seismic stack section of **a** 4°–16°, **b** 15°–25°, and **c** 24°–36°

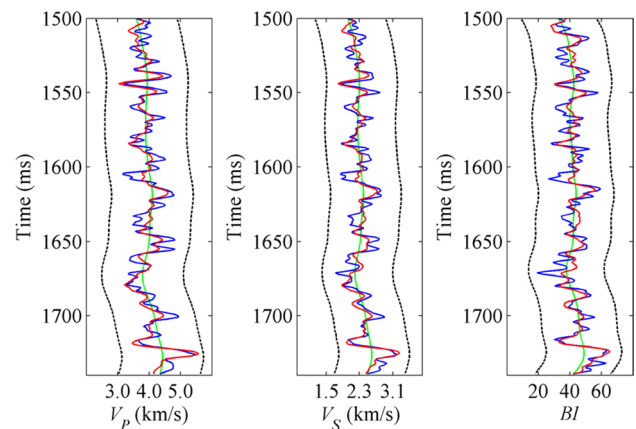
In addition, due to the influence of the population size and the iteration number, the computation cost of NSGA II is very large. Therefore, we only use 1D model, 1D actual logging data and 2D actual seismic data to test the effectiveness of the proposed method on a personal computer with Intel Core CPU i7-8700, NVidia GeForce GTX 1050TI and 8 GB of RAM. However, this method can also be applied to perform the inversion of three-dimensional actual seismic data to estimate BI, which needs further research on high-performance clusters.

### Conclusion

We propose an AVO inversion method for BI based on BI\_Zoeppritz equation and NSGA II. In order to directly estimate the BI, we derive the exact Zoeppritz equation by expressing the density as a function of  $V_p$ ,  $V_s$  and BI and obtain the Zoeppritz equation for BI (BI\_Zoeppritz equation). We also give the analytic expression of  $R_{pp}$  about  $V_p$ ,  $V_s$  and BI. This equation does not introduce any hypothesis, so the  $R_{pp}$  can be calculated from  $V_p$ ,  $V_s$  and BI without



**Fig. 24** Inversion results section of **a**  $V_p$ , **b**  $V_s$  and **c** BI



**Fig. 25** Inversion results of **a**  $V_p$ , **b**  $V_s$  and **c** BI for the seismic trace near the borehole (CDP 1528). The green line and blue line denote the initial model and actual logging data, and the black line and red line denote the search windows and corresponding inverted elastic parameters

loss of accuracy in the inversion process. Because the AVO inversion is a nonlinear problem for three elastic parameters,

we consider it as a global optimization problem of multi-objective function and employ NSGA II to minimize multi-objective function at the same time. To reduce the computation time and improve the inversion accuracy, the proposed method can construct an initial model and an optimized search window to constrain the entire inversion process.

The validity and practicability of the proposed method are verified by the test of theoretical model and actual data. According to the test results, it can be found that the search window has a big impact on inversion accuracy, and the optimized search window can get better inversion results than the linear search window. The appropriate number of genetic iterations can significantly reduce the computation cost while ensuring convergence and inversion accuracy. In addition, the inversion results of angle gathers with  $\text{SNR} = 3$  and  $\text{SNR} = 1$  show that the proposed method has good noise immunity. The application of actual data indicates that it is feasible to obtain BI from actual seismic data using the proposed method.

**Acknowledgements** This work is funded by the SINOPEC's Scientific and Technological Development Program of China (No. P18075-6).

## Compliance with ethical standards

**Conflict of interest** On behalf of all authors, the corresponding author states that there is no conflict of interest.

## References

- Aki K, Richards PG (1980) Quantitative seismology: theory and methods. W H Freeman and Co, Cambridge, pp 144–154
- Deb K, Agrawal RB (1995) Simulated binary crossover for continuous search space. *Complex Syst* 9:115–148
- Deb K, Agrawal S (1999) A niched-penalty approach for constraint handling in genetic algorithms. In: International conference on artificial neural networks and genetic algorithms, pp 235–243
- Deb K, Pratap A, Agarwal S, Meyarivan T (2002) A fast and elitist multi-objective genetic algorithm: NSGA-II. *IEEE Trans Evol Comput* 6(2):182–197
- Fang Y, Zhang FQ, Wang YC (2016) Generalized linear joint PP–PS inversion based on two constraints. *Appl Geophys* 13(1):103–115
- Fu H, Wang X, Zhang L, Gao R, Li Z, Zhu X, Xu W, Li Q, Xu T (2015) Geological controls on artificial fracture networks in continental shale and its fracability evaluation: a case study in the Yanchang Formation, Ordos Basin, China. *J Nat Gas Sci Eng* 26:1285–1293
- Gholami R, Rasouli V, Sarmadivaleh M, Minaeian V, Fakhari N (2016) Brittleness of gas shale reservoirs: a case study from the north Perth basin, Australia. *J Nat Gas Sci Eng* 33:1259–1277
- Glorioso JC, Rattia A (2012) Unconventional reservoirs: basic petrophysical concepts for shale gas. In: SPE/EAGE European unconventional resources conference & exhibition—from potential to production
- Goldberg DE (1989) Genetic algorithms in search, optimization and machine learning. Addison-Wesley, Boston
- Grigg M (2004) Emphasis on mineralogy and basin stress for gas shale exploration. In: SPE meeting on gas shale technology exchange
- Guo TL, Zhang HR (2014) Formation and enrichment mode of Jiaoshiha shale gasfield, Sichuan Basin. *Petrol Explor Dev* 41(1):28–37
- Jarvie DM, Hill RJ, Ruble TE, Pollastro RM (2007) Unconventional shale-gas systems: the Mississippian Barnett Shale of north-central Texas as one model for thermogenic shale-gas assessment. *AAPG Bull* 91(4):475–499
- Jia CZ (2017) Breakthrough and significance of unconventional oil and gas to classical petroleum geological theory. *Petrol Explor Dev* 44(1):1–10
- Jin X, Shah SN, Roegiers JC, Zhang B (2014a) Fracability evaluation in shale reservoirs—an integrated petrophysics and geomechanics approach. In: Proceedings of the SPE hydraulic fracturing technology conference
- Jin X, Shah SN, Truax JA, Roegiers JC (2014b) A practical petrophysical approach for brittleness prediction from porosity and sonic logging in shale reservoirs. In: SPE annual technical conference and exhibition
- Li T, Mallick S (2015) Multicomponent, multi-azimuth prestack seismic waveform inversion for azimuthally anisotropic media using a parallel and computationally efficient non-dominated sorting genetic algorithm. *Geophys J Int* 200(2):1134–1152
- Li JN, Wang SX, Dong CH, Yuan SY, Wang JB (2016) Study on frequency-dependent characteristics of spherical-wave PP reflection coefficient. *Chin J Geophys* 59(10):3810–3819
- Liu W, Wang YC (2018) Multicomponent prestack joint AVO inversion based on exact Zoeppritz equation. *J Appl Geophys* 159:69–82. <https://doi.org/10.1016/j.jappgeo.2018.07.017>
- Mcglade C, Speirs J, Sorrell S (2013) Unconventional gas—a review of regional and global resource estimates. *Energy* 55(1):571–584
- Ostrander WJ (1984) Plane wave reflection coefficients for gas sands at non-normal angles of incidence. *Geophysics* 49(10):1637–1648
- Pei P, Ling K, Hou X, Nordeng S, Johnson S (2016) Brittleness investigation of producing units in Three Forks and Bakken formations, Williston basin. *J Nat Gas Sci Eng* 32:512–520
- Rickman R, Mullen MJ, Petre JE, Grieser WV, Kundert D (2008) A practical use of shale petrophysics for stimulation design optimization: all shale plays are not clones of the Barnett Shale. In: Proceedings of the SPE annual technical conference and exhibition. Society of Petroleum Engineers
- Virieux J, Operto S (2009) An overview of full waveform inversion in exploration geophysics. *Geophysics* 74(6):WCC1–WCC26
- Wang FP, Gale JF (2009) Screening criteria for shale-gas systems. *Gulf Coast Assoc Geol Soc Trans* 59:779–793
- Xie W, Wang YC, Liu XQ, Bi CC, Zhang FQ, Fang Y, Tahir A (2019) Nonlinear joint PP–PS AVO inversion based on improved Bayesian inference and LSSVM. *Appl Geophys* 16(1):64–76
- Yin XY, Liu XJ, Zong ZY (2015) Pre-stack basis pursuit seismic inversion for brittleness of shale. *Petrol Sci* 12(4):618–627
- Yuan SY, Liu Y, Zhang Z, Luo CM, Wang SX (2019) Prestack stochastic frequency-dependent velocity inversion with rock-physics constraints and statistical associated hydrocarbon attributes. *IEEE Geosci Remote Sens Lett* 16(1):140–144
- Zhang D, Ranjith PG, Perera MSA (2016) The brittleness indices used in rock mechanics and their application in shale hydraulic fracturing: a review. *J Petrol Sci Eng* 143:158–170
- Zoeppritz K (1919) On the reflection and propagation of seismic waves. *Gott Nachr* I:66–84
- Zong ZY, Yin XY, Wu GC (2013) Elastic impedance parameterization and inversion with Young's modulus and Poisson's ratio. *Geophysics* 78(6):N35–N42
- Zou CN, Yang Z, Zhang GS, Hou LH, Zhu RK, Tao SZ, Yuan XJ, Dong DZ, Wang YM, Guo QL, Wang L, Bi HB, Li DH, Wu N (2014) Conventional and unconventional petroleum “orderly accumulation”: concept and practical significance. *Petrol Explor Dev* 41(1):14–30



# Interpretation of gravity anomaly over 2D vertical and horizontal thin sheet with finite length and width

Arkoprovo Biswas<sup>1</sup>

Received: 6 May 2020 / Accepted: 7 July 2020 / Published online: 16 July 2020  
© Institute of Geophysics, Polish Academy of Sciences & Polish Academy of Sciences 2020

## Abstract

Gravity data are often used for delineation of the lateral and vertical extension of mineralized bodies buried at different depths. Various parameters associated with the buried bodies are the primary concern for mineral exploration purposes. Hence, a reliable and efficacious interpretation method is developed for the delineation of gravity anomaly data over the 2D vertical and horizontal sheet with finite length and width associated with mineralized bodies. The parameters viz. amplitude coefficient ( $k$ ), location ( $x_0$ ), depth to the top of the body ( $h$ ), length of the sheet ( $L$ ), and shape factor ( $q$ ) for 2D vertical sheet-type structure and depth ( $h$ ) and width ( $w$ ) of the sheet for 2D horizontal sheet were resolved. Restricting  $x_0$  and  $q$  has given very reliable results for the 2D vertical sheet, and the  $w$  for 2D horizontal sheet shows the problem of equivalence. However, in all cases, the delineated parameters are within the expected uncertainty. The present interpretation method was applied to synthetic and noisy data and three field examples from the USA, Canada, and Sweden for mineral exploration purposes. It has also been seen that the present study is more reliable in delineating the actual structure associated with mineralized bodies for the 2D vertical and horizontal sheet-type structure. The delineated parameters are in outstanding agreement with the earlier works, borehole information and also updated the actual subsurface structure.

**Keywords** Gravity anomaly · 2D sheet · VFSA · Mineral exploration

## Introduction

Gravity data have been used for many purposes from crustal studies to exploration of oil, gas, and mineralized bodies. Most of the structures associated with such structures were approximated by perfect geological bodies viz. a sphere, horizontal or vertical cylinder with semi-infinite length and width, dyke and also by a 2D finite sheet-like structures. For such types of structures, numerous interpretation methodologies were developed and the best interpretation for different parameters was delineated.

The gravity anomaly over 2D vertical and horizontal sheet with finite length and width has considerable attention for the interpretation of mineralized bodies (Abdelrahman et al. 2016; Kara and Hoskan 2016; Essa and Geraud 2020) Tlas and Asfahani 2018). Approximation of the depth,

and the exact length and width of any mineralized body is very essential for exploration purposes. Many interpretation approaches were developed for the interpretation of gravity data which can be categorized in three different ways. The first category is the continuous modeling where you need the density distribution along with other geophysical data knowledge to constrain the subsurface structures (Talwani et al. 1959; Holstein et al. 2010; Hinze et al. 2013). The second category was defined earlier by Nabighian (1972, 1974) as either automatic or semi-automatic methods which were also applied by Blakely and Simpson (1986), Reid et al. (1990), Marson and Klingele (1993), Klingele et al. (1991) and Mackleod et al. (1993). The third category is the quantitative interpretation which the present work is based upon is the interpretation of various parameters of idealized subsurface structures (Essa and Munsch 2019; Anderson et al. 2020; Geldart et al. 1966; Green 1976).

Many interpretations were developed in the past considering various approaches starting from curve matching, graphical methods, nomograms, least-square methods, etc. The detailed information can be seen from various literature published earlier (Biswas 2015, 2016; Abdelrahman et al. 2016;

✉ Arkoprovo Biswas  
arkoprovo@gmail.com; arkoprovo.geo@bhu.ac.in

<sup>1</sup> Department of Geology, Centre of Advanced Study, Institute of Science, Banaras Hindu University, Varanasi, UP 221005, India

Tlas and Asfahani 2018, and references therein). Moreover, advanced techniques such as Fourier transform, Euler deconvolution, Mellin transform, neural network, least-squares, minimization approaches, Werner deconvolution, etc., (Odegard and Berg 1965; Sharma and Geldart 1968; Hartmann et al. 1971; Jain 1976; Thompson 1982; Gupta 1983; Lines and Treitel 1984; Mohan et al. 1986; Abdelrahman 1990; Abdelrahman et al. 1991; Abdelrahman and El-Araby 1993; Abdelrahman and Sharafeldin 1995; (Khalil et al. 2015; Khalil et al. 2014; Abdelrahman and Essa 2015; Essa 2014; Abdelrahman and Essa 2013; Abdelrahman et al. 2003 Elawadi et al. 2001) were also developed for the understanding of gravity anomalies caused due to different subsurface structures.

Earlier many subsurface structures interpreted using different approaches considered different structures for the same field anomaly data. This is one of the major problems in delineating the actual subsurface structures associated with a mineralized body. Although the main objective is to find the accurate depth of the body, however, many fail to determine the near probable true structure of the bodies. Gravity data interpretation are always ill-posed, and it is very difficult to determine the actual solution, and sometimes it will give an equivalent solution that could be erroneous (Mehanee 2014; Mehanee and Essa 2015). To overcome these limitations, global optimization methods such as neural network modeling (Abedi et al. 2010), very fast simulated annealing (Biswas 2015, 2016); differential evolution algorithm (Ekinci et al. 2016, 2019; Balkaya et al. 2017), particle swarm optimization (Singh and Biswas 2016; Essa and Munsch 2019; Anderson et al. 2020), ant colony optimization (Srivastava et al. 2014) were applied to interpret gravity data. A review on Global optimization in potential field anomalies is described by Ekinci et al. (2020) and references therein. Hence, in the present work, interpretation of all parameters associated with a 2D vertical and horizontal sheet with finite length and width was carried out. Very fast simulated annealing after Biswas 2015) was applied for the effective interpretation of the gravity anomaly over 2D sheet-type structure besides the uncertainty related with the interpretation of model parameters which was not discussed in earlier literature. The method is confirmed and appraised on synthetic and noisy data and also from three field data from the USA, Canada, and Sweden.

## Methodology

### Formulation of the forward problem

The gravity anomaly over a 2D vertical and horizontal sheet of finite length and width on a plane is measured by the following equation.

### 2D vertical sheet with finite length

Following Nettleton (1942), and Kara and Hoskan (2016), the gravity anomaly produced by a vertical sheet/line of finite length is given by (Fig. 1a)

$$g(x_i) = k \left[ \frac{1}{\{(x_i - x_0)^2 + h^2\}^{1/2}} - \frac{1}{\{(x_i - x_0)^2 + (h + L)^2\}^{1/2}} \right] \quad (1)$$

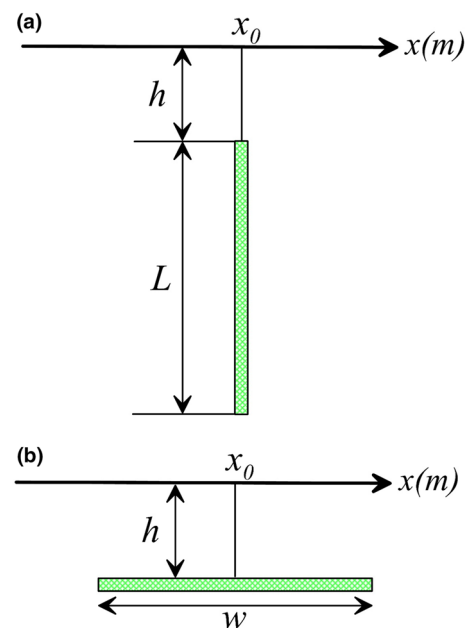
$i = 1, 2, 3, \dots, N$

where  $k$  is the amplitude coefficient ( $k = GP$ , where  $G$  is gravity constant,  $\rho$  is the linear density of anomalous mass in the line),  $h$  is the depth from the top of the sheet,  $L$  is the length of the sheet, and  $x_0$  is the origin of the sheet on the surface. The above Eq. 1 can also be written as

$$g(x_i) = k \left[ \frac{1}{\{(x_i - x_0)^2 + h^2\}^q} - \frac{1}{\{(x_i - x_0)^2 + (h + L)^2\}^q} \right] \quad (2)$$

$i = 1, 2, 3, \dots, N$

where other parameters remain the same except the power (1/2) in Eq. 1 is taken as another parameter ( $q$ ), where  $q=0.5$  and is a dimensionless quantity and can also be taken as a shape factor. All five parameters will be interpreted in the present work.



**Fig. 1** A 2D sheet type structure within the subsurface **a** vertical with finite length, and **b** horizontal with finite width

## 2D horizontal sheet with finite width

Following Pick et al. (1973), and Abdelrahman et al. (2016), the general expression for gravity anomaly produced over a horizontal sheet/line of finite width is given by (Fig. 1b)

$$g(x_i) = k \left[ \tan^{-1} \left( \frac{w - 2(x_i - x_0)}{2h} \right) + \tan^{-1} \left( \frac{w + 2(x_i - x_0)}{2h} \right) \right]$$

$$i = 1, 2, 3, \dots, N \quad (3)$$

where  $k$  is the amplitude coefficient ( $k = 2G\sigma t$ , where  $G$  is the gravitational constant, and  $\sigma t$  is the surface density contrast of the sheet where  $t$  is the thickness of the thin sheet),  $h$  is the depth from top to the middle point of the sheet,  $w$  is the width of the sheet, and  $x_0$  is the origin of the sheet on the surface. However, this Eq. 3 is only effective once the thickness of the sheet is very insignificant as compared to its width.

## Inversion using global optimization

The global optimization method was developed to find out the best model (Global solution) in the presence of several local optima. It is a controlled optimization approach that tries to find out the globally best solutions. It has been applied in many applications and since geophysical models are mostly nonlinear, and hence, it needs a global solution. In the present work, an advanced version of simulated annealing (SA) termed very fast simulated annealing (VFSA) will be used in the gravity data.

### Very fast simulated annealing (VFSA)

SA or VFSA was resulting from the similarity of the heat bath algorithm (Sen and Stoffa 2013) and was used for the delineation of much geophysical data. An overview of VFSA is described in many published literature (Sen and Stoffa 2013; Sharma 2012; Biswas 2015) and henceforth not described here for brevity. VFSA was developed to negate the effect of linear inversion approach since all optimization problems are not linear. The method is very robust and has high stability and enhanced resolution (in the presence of multiple local optima) of the solution. It can also negate the problem of non-uniqueness and takes very less computing time in extremely large data, its resolution of the data and can find the global model (solution/s) (Sen and Stoffa 2013). Also, Ingber and Rosen (1992) mentioned that VFSA takes very little CPU time and memory and gives high-resolution results.

Here, the calculation of misfit error for delineations of gravity anomaly is taken as (after Sharma and Biswas 2013)

$$\varphi = \frac{1}{N} \sum_{i=1}^N \left( \frac{D_i^0 - D_i^c}{|D_i^0| + (D_{\max}^0 - D_{\min}^0)/2} \right)^2 \quad (4)$$

where  $N$  refers to the number of data,  $D_i^0$  and  $D_i^c$  is the  $i$ th observed data and model responses (data),  $D_{\max}^0$  and  $D_{\min}^0$  are the maximum (+) and minimum (−) value of the synthetic/field data.

## Global model and uncertainty analysis

Global model/solutions and uncertainty analysis were applied in various fields of study and also in the interpretation of nonlinear geophysical data (Ekinici et al. 2019, 2020; Fernández-Martínez et al. 2020). So, a globally best model and uncertainty analysis are a must for every nonlinear geophysical inversion. Hence, the procedure developed by Mosegaard and Tarantola (1995) and Sen and Stoffa (1996) was applied in the present study. Moreover, uncertainty analysis and the probability density function (PDF) for every model parameter were studied following Trivedi et al. (2020). The details of the global model and uncertainty analysis are described in various kinds of literature (Sharma 2012; Sharma and Biswas 2013) and not conversed here for conciseness. The developed algorithm for interpretation of gravity anomaly data was performed in a Windows 10 platform by MS FORTRAN Developer studio through the Intel Core i7 processor with a CPU time of 45 s.

## Results and discussion

### Synthetic example

#### Parameter search and tuning

In the interpretation of gravity or other potential field data caused by some idealized buried structures, it is known that the location of the body ( $x_0$ ) could be understood from the highest/lowest peak, depth of the body ( $h$ ) from the surface can be predicted from the half-width, and the amplitude coefficient ( $k$ ) from the highest or lowest value (Nabighian 1972) for the qualitative delineation. This was shown in various literature published earlier; however, Srivastava and Agarwal (2010) observed that the horizontal location ( $x_0$ ) and the shape factor ( $q$ ) are the most stable parameters. This was also shown in the interpretation of gravity anomaly data (Biswas 2015, 2016; Biswas et al. 2017; Trivedi et al. 2020). Hence, in the present work, a two-step process was executed. To delineate the best interpretation for all the parameters, the search space was kept in a large range. After a single



run of the inversion, if the interpreted parameters are within the specified range, the search space was condensed and the inversion procedure was again repeated (10 runs). Next, it was observed as mentioned earlier that the interpretation of parameter ( $x_0$ ) is the same as the original value is chosen and ( $q$ ) is very close to the same value taken. Hence, in the next step, both the parameters ( $x_0$  and  $q$ ) were restricted to its original value and the inversion process was carried out. This procedure shows that all the other parameters ( $k$ ,  $h$ , and  $L$ ) were interpreted accurately with less error and uncertainty. This procedure was followed for every synthetic model and also for field data associated with vertical 2D sheet with finite length. In the case of a 2D sheet with finite width, all the parameters were interpreted and the results were discussed. Here, for 2D finite width sheet,  $x_0$  also showing the same as discussed for 2D finite length sheet. However, other parameters were interpreted without restricting  $x_0$  to its original value.

**Model 1 (vertical sheet)**

A 2D inclined sheet with a finite length model was taken, and the anomaly (Fig. 2) was generated using Eq. 1. The two-step inversion process was executed, and the parameters ( $k$ ,  $x_0$ ,  $h$ ,  $L$ ,  $q$ ) were delineated (Table 1). Histograms were also prepared for all parameters and are shown in Fig. 3a. It can be seen that the parameters  $x_0$  and  $q$  are very close to the true value. Hence, in the next step, the parameters  $x_0$  and  $q$  were restricted to its original value and the inversion process was again carried out. After restricting the parameters  $x_0$  and  $q$ , the histogram was again prepared and now it can be seen that the other parameters ( $k$ ,  $h$ , and  $L$ ) are near to its original value (Fig. 3b).

Moreover, to see the effect of noise, 10% Gaussian noise was added in Model 1 and the inversion process was repeated. Figure 3c, d shows the histograms from all parameters and after restricting  $x_0$  and  $q$  to its original value. From the histogram analysis, it can be concluded that the two-step inversion process can exactly delineate all the model

parameters. Responses from synthetic data without noise and noisy data and calculated models are seen from Fig. 2a, b. The final interpreted parameters with errors are shown in Table 1, and it can be seen that following this two-step procedure, the estimated errors can be reduced and the parameters are almost the true value.

A 3D cross-plot study was also accomplished to see the relationship between different parameters and their

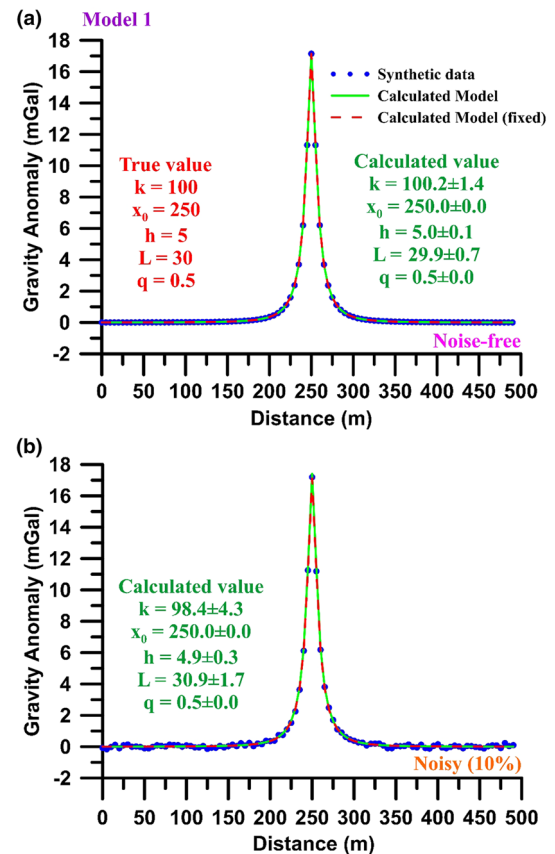


Fig. 2 Calculated model response for Model 1 (vertical sheet)

**Table 1** Inversion results from Model 1

Parameters	True model	Search space	Synthetic data	Synthetic data (controlled $x_0$ and $q$ )	Noisy data	Noisy data (controlled $x_0$ and $q$ )
$k$ (mGal)	100	0–200	$80.9 \pm 7.5$	$100.2 \pm 1.4$	$116.9 \pm 17.8$	$98.4 \pm 4.3$
$x_0$ (m)	250	0–300	$250.0 \pm 0.0$	$250.0 \pm 0.0$	$250.0 \pm 0.1$	$250.0 \pm 0.0$
$h$ (m)	5	0–10	$4.7 \pm 0.2$	$5.0 \pm 0.1$	$5.2 \pm 0.3$	$4.9 \pm 0.3$
$L$ (m)	30	0–50	$25.6 \pm 1.9$	$29.9 \pm 0.7$	$34.3 \pm 3.8$	$30.9 \pm 1.7$
$q$	0.5	0–2	$0.42 \pm 0.0$	$0.5 \pm 0.0$	$0.55 \pm 0.1$	$0.5 \pm 0.0$
Error			$1.8 \times 10^{-6}$	$3.7 \times 10^{-9}$	$9.0 \times 10^{-5}$	$8.3 \times 10^{-5}$

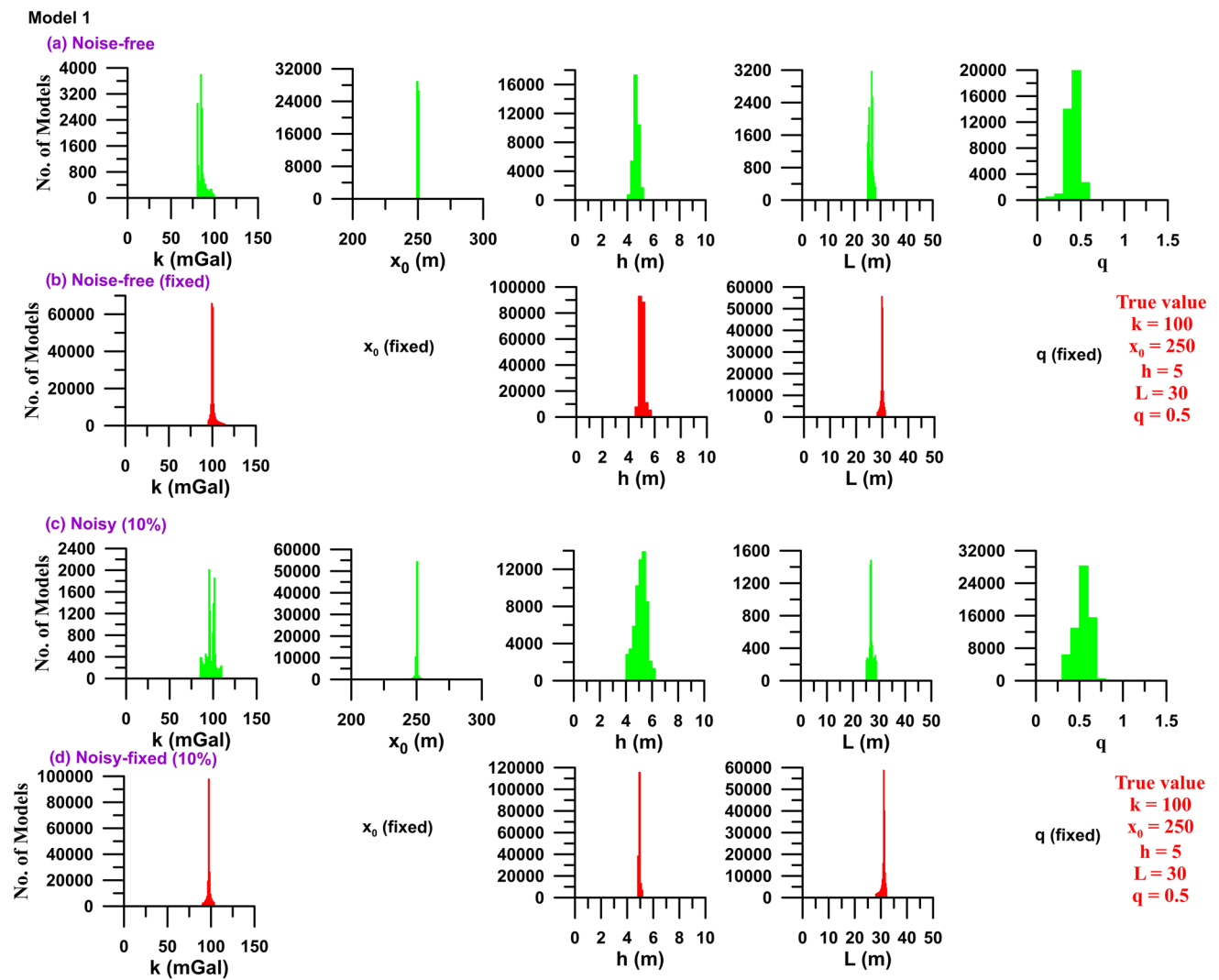


Fig. 3 Histogram study for Model 1 (vertical sheet)

uncertainty associated with the final interpretation. It has been seen from the histogram study that restricting  $x_0$  and  $q$  will give better results. Hence, initially, when all the parameters were interpreted, a cross-plot between  $k$ ,  $h$ , and  $L$  was prepared. It can be seen from Fig. 4a that the model shows a high range. In the next step, the parameters  $x_0$  and  $q$  were restricted to its true value and again the cross-plots between  $k$ ,  $h$ , and  $L$  were made (Fig. 4b). It can be seen that after restricting  $x_0$  and  $q$  to its true value, the parameters are very near to its original value and the uncertainty decreases. This can also be seen from the results as shown in Table 1. Furthermore, to check whether the existence of noise impacts the parameters, noise corrupted data were also evaluated and

the outcomes were very near to the true value. Figure 4c, d shows the cross-plots for noisy data. Cross-plots study advocates that the model parameters are very adjacent to the original value (Yellow) and the ultimate mean model parameters remained in the uncertainty limitations (one standard deviation) and also in the peak PDF (Red).

**Model 2 (vertical sheet)**

An alternate model (Table 2) was taken with a high amplitude coefficient, more depth, and a larger length of the sheet, and the anomaly was generated using Eq. 2. The two-step inversion process was also executed for this model (Fig. 5a).

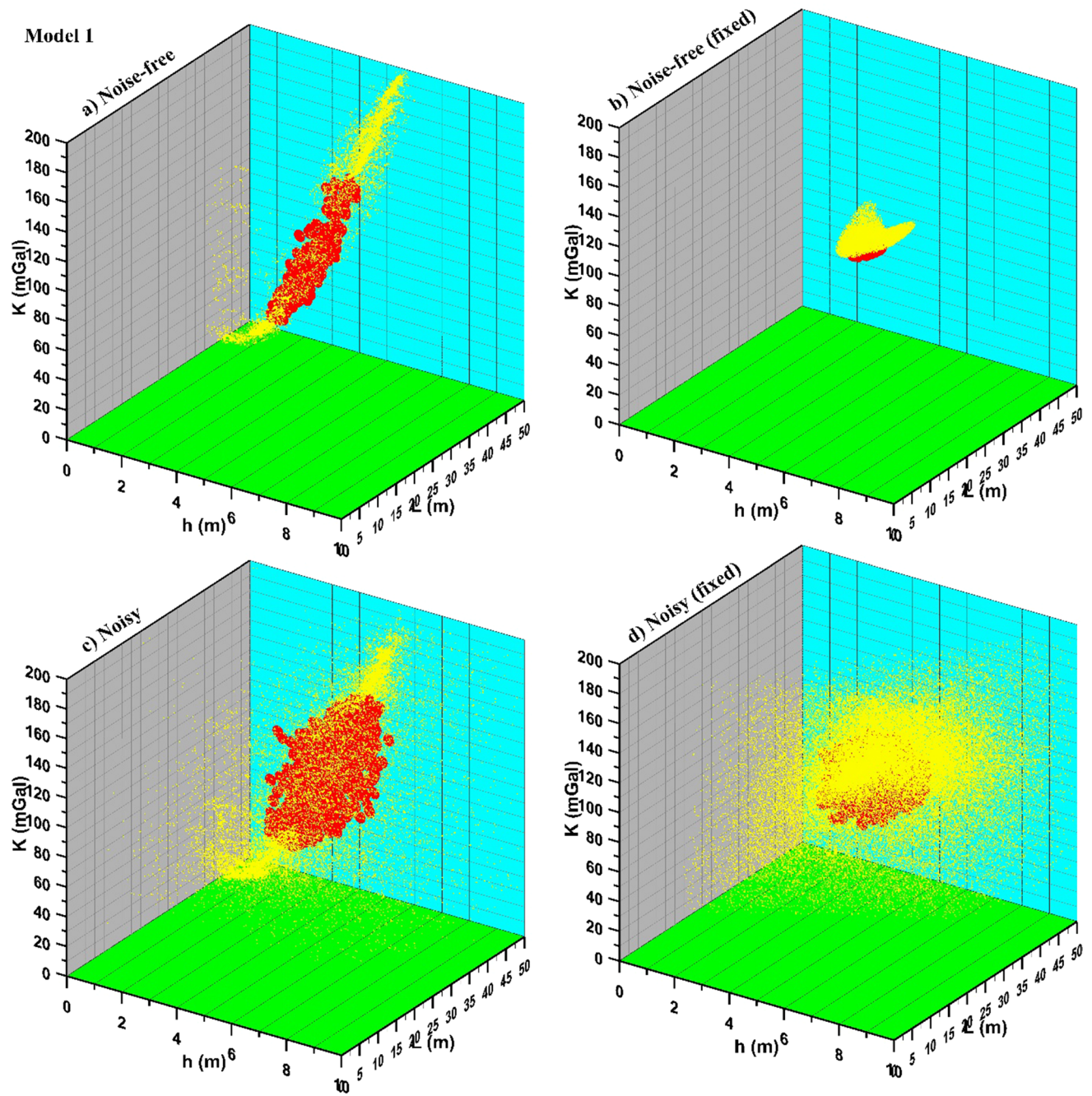
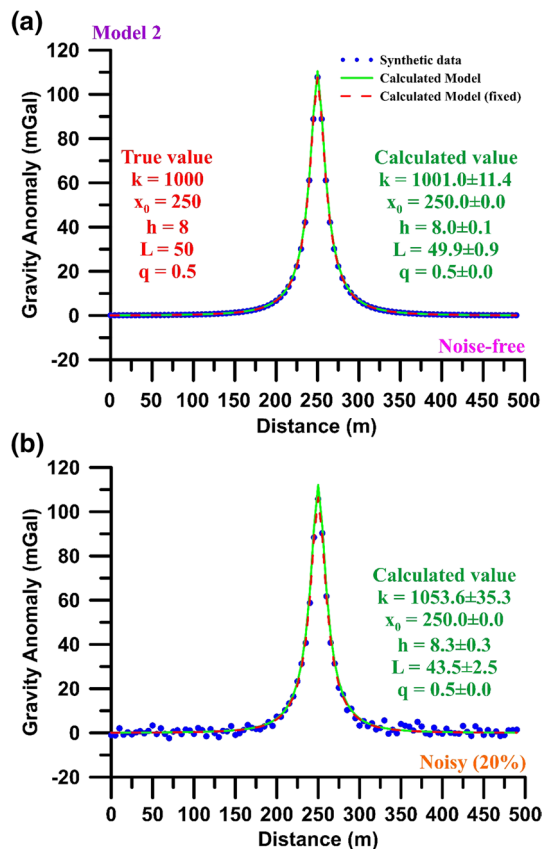


Fig. 4: 3D cross-plot for Model 1 (vertical sheet)

**Table 2** Inversion results from Model 2

Parameters	True model	Search space	Synthetic data	Synthetic data (controlled $x_0$ and $q$ )	Noisy data	Noisy data (controlled $x_0$ and $q$ )
$k$ (mGal)	1000	0–1500	$851.2 \pm 142.99$	$1001.0 \pm 11.4$	$996.3 \pm 224.8$	$1053.6 \pm 35.3$
$x_0$ (m)	250	0–300	$250.0 \pm 0.0$	$250.0 \pm 0.0$	$250.0 \pm 0.2$	$250.0 \pm 0.0$
$h$ (m)	8	0–10	$7.7 \pm 0.3$	$8.0 \pm 0.1$	$8.2 \pm 0.5$	$8.3 \pm 0.3$
$L$ (m)	50	0–100	$45.5 \pm 4.5$	$49.9 \pm 0.9$	$41.2 \pm 6.3$	$43.5 \pm 2.5$
$q$	0.5	0–2	$0.45 \pm 0.1$	$0.5 \pm 0.0$	$0.47 \pm 0.1$	$0.5 \pm 0.0$
Error			$2.0 \times 10^{-5}$	$1.8 \times 10^{-9}$	$7.3 \times 10^{-4}$	$6.3 \times 10^{-4}$

Moreover, to check the effect of higher degrees of noises, 20% Gaussian noise was also added in the data and the inversion procedure was also carried out (Fig. 5b). Histogram analysis was also carried out for this model, and results reveal the same (Fig. 6a, b) as discussed for Model 1. Figure 6c, d shows the histogram derived from noisy data and also concludes that all the parameters were delineated correctly. The final interpreted models are shown in Table 2, and the responses from synthetic and noisy data are shown in Fig. 5a, b.

**Fig. 5** Calculated model response for Model 2 (vertical sheet)

### Model 3 (horizontal sheet)

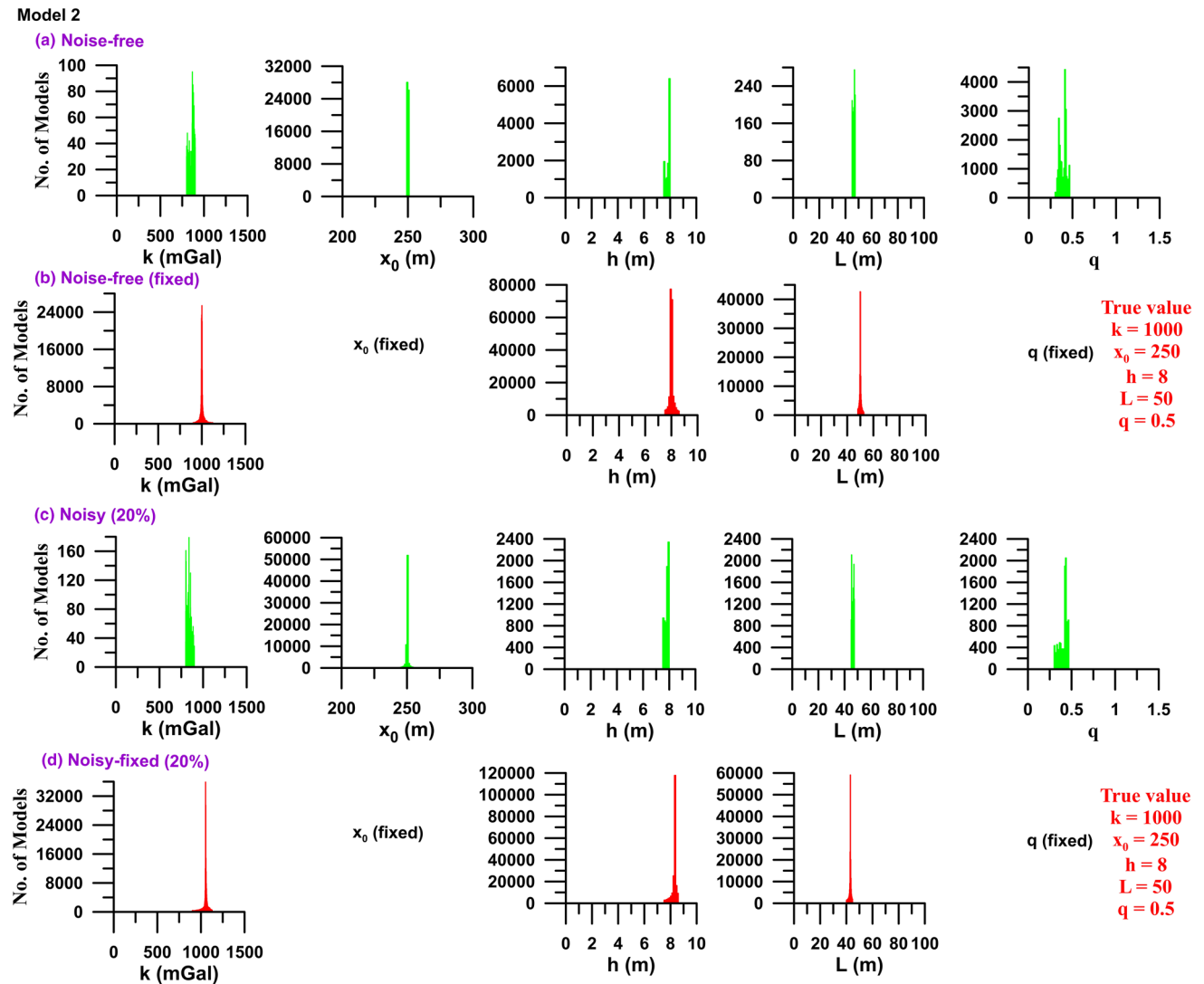
This model is taken for a 2D horizontal sheet with finite width, and the anomaly was generated using Eq. 3 (Fig. 7). The inversion procedure was carried out for this model, and all the parameters were delineated. Histogram for this model was also prepared, and it can be seen that the parameters  $x_0$  and  $h$  were resolved very efficiently. However, there is a bit large range for  $k$  and  $w$  (Fig. 8a). This suggests that the location and the depth of the body are well solved; however, there is some uncertainty in delineating the other two parameters. Moreover, the study of the histogram for noisy data also tells the same as discussed (Fig. 8b).

Moreover, to study the uncertainty and relationship amongst different parameters, a 3D cross-plot is also studied for synthetic data without noise (Fig. 9a). From the cross-plot, it can be seen that  $h$  is well resolved, but for the specific value of  $h$ , there is a large range of  $w$  and  $k$ . This also suggests that for a specific depth ( $h$ ) of the body, the width ( $w$ ) of the body can change, and it will fit the data exactly well. This confirms that there is an equivalence problem associated with such type of structure. Moreover, the cross-plot study was also studied for noisy data and it also shows the same as discussed above (Fig. 9b). Cross-plots study demonstrates that parameters ( $k$ ,  $h$ , and  $w$ ) are close to the initial value (Yellow). Mean model parameters endured in the uncertainty bounds (one standard deviation) and also in the highest PDF (Red). The concluding interpreted parameters for this model are shown in Table 3, and the calculated models from both the data are shown in Fig. 7a, b.

### Field example

#### Louga anomaly, USA

The field example is taken from the Louga gravity anomaly, USA (after Nettleton 1976, Fig. 14.8 therein). This is a north–south gravity profile. The data were digitized from



**Fig. 6** Histogram study for Model 2 (vertical sheet)

Nettleton (1976) at the same interval. The data were interpreted using the present inversion method. It has been found that the anomaly is due to a vertical 2D sheet. The depth and length of the sheet interpreted to be 5 and 37.2 km, respectively. The field data were earlier interpreted as a spherical structure by Mohan et al. (1986) using Mellin transformation and a finite line by Kara and Hoskan (2016). The depth obtained from different methods is shown in Table 4. The comparison between field data and calculated model as well as the subsurface structure is shown in Fig. 10 and compared with Kara and Hoskan (2016).

### Mobrun anomaly, Noranda, Quebec, Canada

The gravity anomaly example was taken over a massive sulfide ore body from Noranda Mining District, Quebec, Canada (Siegel et al. 1957; Grant and West 1965). The field data were taken from Grant and West (1965) and digitized at the same interval. The field data were earlier interpreted by several workers using diverse interpretation approaches (Skeels 1963; Roy 1966; Atchuta Rao et al. 1985; Sundararajan et al. 2000). However, the same field data were also interpreted lately by Biswas (2015) considering horizontal

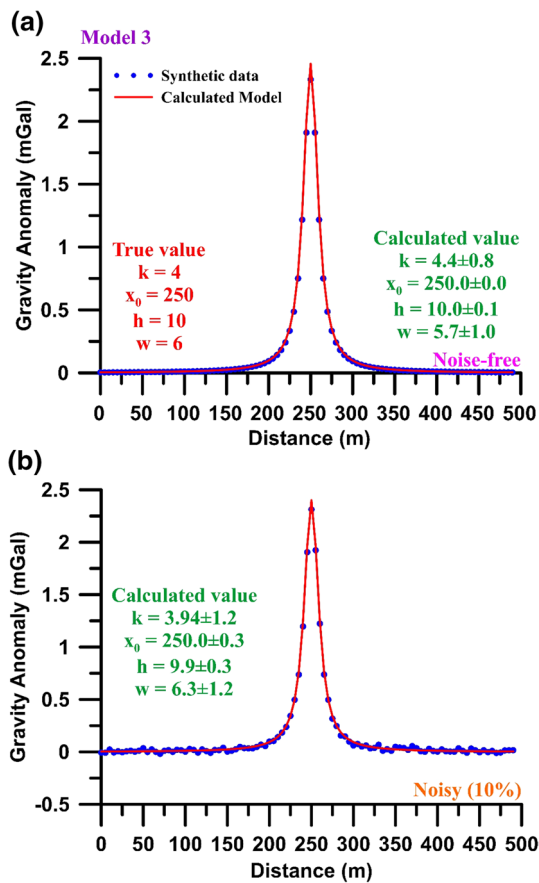


Fig. 7 Calculated model response for Model 3 (horizontal sheet)

cylinder, Biswas (2016) considering an inclined infinite sheet using VFSA optimization, Singh and Biswas (2016) using particle swarm optimization, and Al-Garni (2018) considering inclined finite sheet using integral transform method. The data were again interpreted considering a 2D vertical finite sheet in the present case, and it was found that the depth obtained from the present study is 31.8 m and the length of the sheet is 161.4 m. The interpreted depth obtained in the present study is very close to the drilling data (30.48 m). The field data were also interpreted by Siegel et al. (1957) and Skeels (1963) considered a vertical prism. They have also estimated the depth to the bottom of the ore bodies. The interpreted depth and the length of the body are recalculated in the present work and shown in Table 5. It must be mentioned that the ore body present in that area is vertical (see Fig. 9; Roy 1966). It should also be mentioned that earlier the same anomaly was interpreted considering infinite horizontal cylinder, inclined infinite sheet. However, the results obtained in those works are conclusive, but the same anomaly also suggests that it has a depth extension and the length of the ore body can also be determined. The error estimated from the present study is less compared to other interpretation methods considering different subsurface structures. A comparison of different estimated parameters is shown in Table 5. Figure 11 shows the comparison of field data and model responses

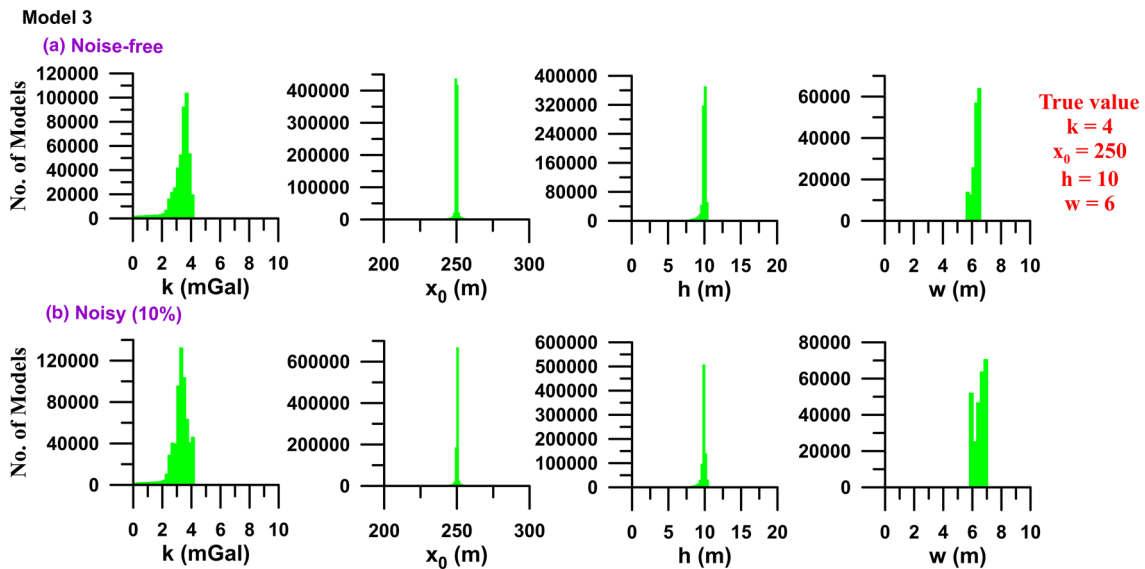


Fig. 8 Histogram study for Model 3 (horizontal sheet)

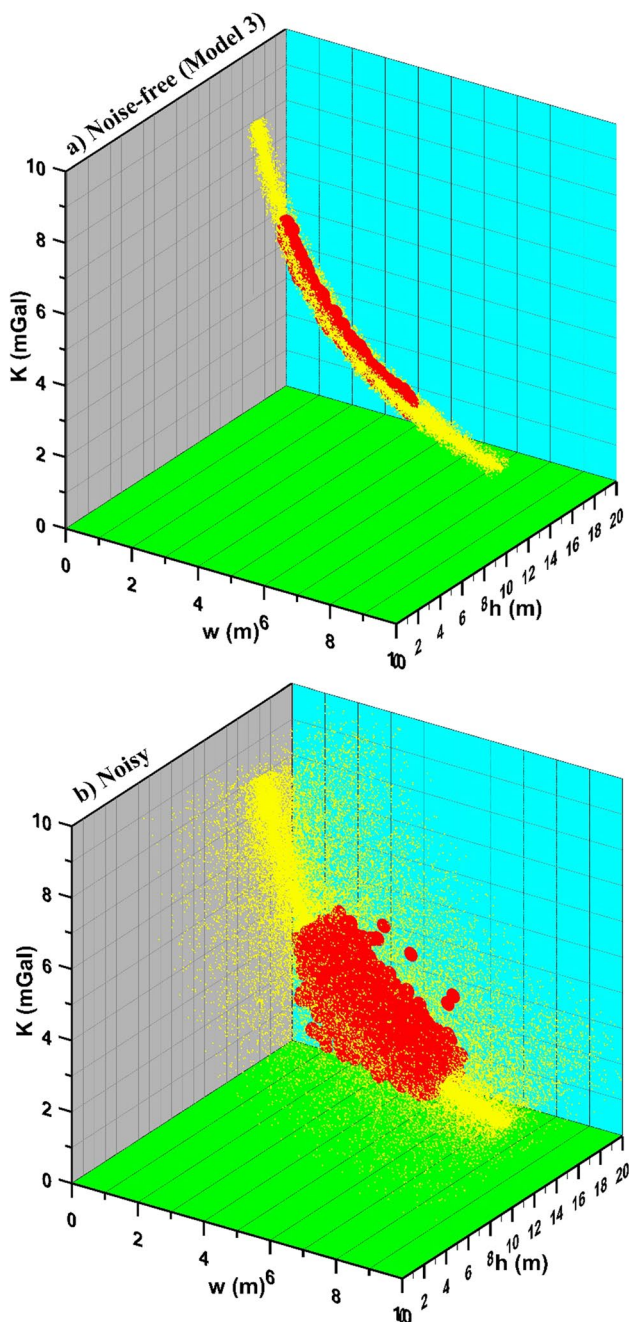


Fig. 9 3D cross-plot for Model 3 (horizontal sheet)

**Table 3** Inversion results from Model 3

Parameters	True model	Search space	Synthetic data	Noisy data
$k$ (mGal)	4	0–10	$4.4 \pm 0.8$	$3.9 \pm 1.2$
$x_0$ (m)	250	0–300	$250.0 \pm 0.0$	$250.0 \pm 0.3$
$h$ (m)	10	0–20	$10.0 \pm 0.1$	$9.9 \pm 0.3$
$w$ (m)	6	0–10	$5.7 \pm 1.0$	$6.3 \pm 1.2$
Error			$1.5 \times 10^{-6}$	$1.5 \times 10^{-3}$

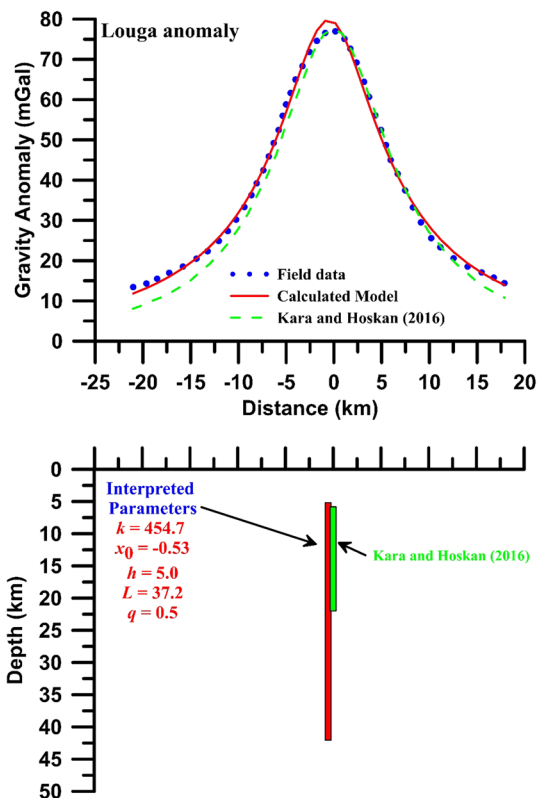
from present and earlier works along with the interpreted subsurface structure.

### Karrbo gravity anomaly, Sweden

The last gravity anomaly data over the pyrrhotite ore body at Karrbo, Vastmanland, Sweden (Shaw and Agarwal 1990) were taken for interpretation (Fig. 12). The data were taken from Shaw and Agarwal (1990) and digitized for the same interval. The data were earlier interpreted as horizontal cylinder/infinite horizontal rod (Tlas and Asfahani 2018; Singh and Biswas 2016; Biswas 2015; Asfahani and Tlas 2012; Tlas et al. 2005); however, the width of the body was never interpreted. Hence, the field data were tried to interpret considering a horizontal sheet with a finite width. The interpretation of the data shows that the depth of the body is 4.6 m and with is 1.8 m, respectively. It can be well understood from Table 6 that the depth and the location of the ore body are well delineated considering different structures and interpretation methods. However, in earlier cases, the width of the body is not delineated. Moreover, the responses from this interpretation perfectly fit with the field data (Fig. 12) which is also the case with earlier interpretation. Hence, it can be said that although the depth is well resolved in all the earlier interpretations, considering the width of the body it can also be delineated approximately although the width has some uncertainty the information can be used drilling and exploration purposes.

**Table 4** Inversion results from Louga anomaly, USA

Parameters	Search space	Present study—(vertical finite sheet)	Kara and Hoskan (2016)—(vertical finite line)	Mohan et al. (1986)—(sphere)	Nettleton (1976)
$k$ (mGal)	0–1000	$454.7 \pm 10.9$	602	–	–
$x_0$ (km)	–1 to 1	$-0.53 \pm 0.0$	–	–	–
$h$ (km)	0–10	$5.0 \pm 0.1$	5.75	9.31	9.30
$L$ (km)	0–50	$37.2 \pm 1.9$	16.3	–	–
$q$	0–2	$0.5 \pm 0.0$	0.5	–	–
Error		$3.3 \times 10^{-4}$	–	–	–

**Fig. 10** Calculated model response and subsurface structure for Louga anomaly, USA

## Conclusion

Interpretation of gravity anomaly was carried out for 2D vertical and horizontal sheet type structures with a finite length and width of the body. The present inversion method can delineate the amplitude coefficient ( $k$ ), location ( $x_0$ ), depth to the top of the body ( $h$ ), length of the sheet ( $L$ ), and shape factor ( $q$ ) for 2D vertical sheet type structure and depth ( $h$ ) and width ( $w$ ) of the sheet for 2D horizontal sheet. The inversion results from 2D vertical sheet shows that restricting the  $x_0$  and  $q$  to its original value taken would give the most reliable results. In case of 2D horizontal sheet, all the parameters are well delineated; however, the width ( $w$ ) of the sheet shows some uncertainty, i.e., the parameters show an equivalent solution. However, in both the cases for 2D sheet, all the delineated parameters are within the uncertainty limits. Histograms and cross-plots were also studied for both the cases, and it also shows the same. The present interpretation method has been verified with synthetic data with no noise and different degrees of Gaussian noise added in the data. Moreover, three field examples from different locations were interpreted in terms of mineral bodies delineation. The delineated parameters achieved by the current technique can be applied to know the subsurface structure associated with mineral deposits.



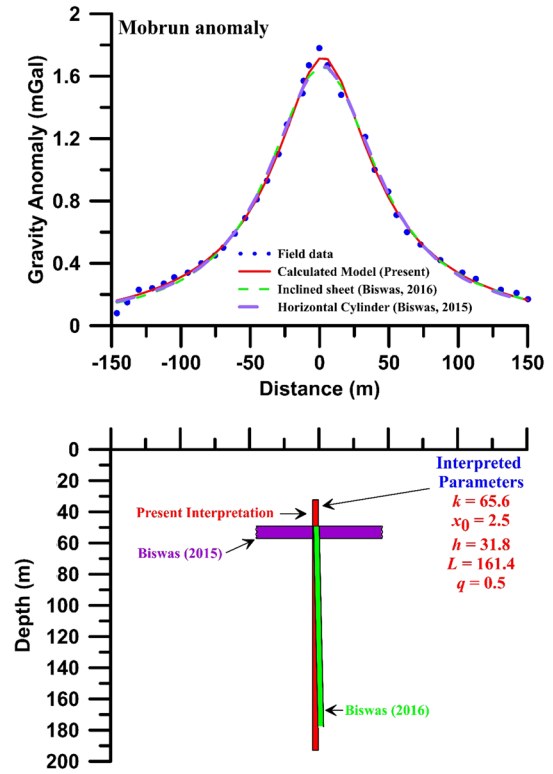
**Table 5** Inversion results from Mobrun Anomaly, Noranda, Quebec, Canada

Parameters	Search space	Present study—(vertical finite sheet)	Al-Garni (2018) (inclined infinite sheet)	Singh and Biswas (2016) (horizontal cylinder)	Biswas (2016) (inclined infinite sheet)	Biswas (2015) (horizontal cylinder)	Mehanee (2014)	Sundara-rajana et al. (2000)	Atchuta Rao et al. (1985)	Roy (1966)	Grant and West (1965)	Skeels (1963)	Siegel et al. (1957)
$k$ (mGal)	0–100	$65.6 \pm 1.5$	146.37	—	$79.5 \pm 0.5$	$79.5 \pm 0.7$	80.0	—	—	—	—	—	—
$x_0$ (m)	0–5	$2.5 \pm 0.2$	—	2.37	$1.4 \pm 0.4$	$2.5 \pm 0.4$	—	—	—	—	—	—	—
$h$ (m)	0–60	$31.8 \pm 0.6$	28.04	46.69	$47.9 \pm 0.4$	$47.7 \pm 0.6$	47.0	18.2	16.25	21.34	17.5	25.6	6.0
$L$ (m)	0–200	$161.4 \pm 5.6$	—	1.0	—	1.0	—	—	—	—	—	$145.0^a$	$192.1^b$
$q$	0–2	$0.5 \pm 0.0$	—	—	—	1.0	1.0	—	—	—	—	—	—
Error		$5.0 \times 10^{-4}$	—	—	$6.2 \times 10^{-4}$	$6.5 \times 10^{-4}$	—	—	—	—	—	—	—

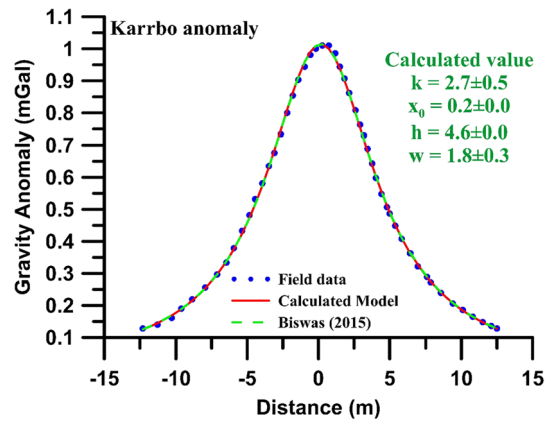
$h = 100$  ft (30.48 m) from drilling data

<sup>a</sup> Depth to the bottom is 170.68 m (length calculated: depth to the bottom – depth to the top)

<sup>b</sup> Depth to the bottom is 198.12 m (length calculated: depth to the bottom – depth to the top)



**Fig. 11** Calculated model response and subsurface structure for Mobrun Anomaly, Noranda, Quebec, Canada



**Fig. 12** Calculated model response for Karrbo Gravity Anomaly, Sweden

**Table 6** Inversion results from Karrbo Gravity Anomaly, Sweden

Parameters	Search space	Present study—(horizontal finite sheet)	Tlas and Asfahani (2018)	Singh and Biswas (2016)	Biswas (2015)	Asfahani and Tlas (2012)	Tlas et al. (2005)
$k$ (mGal)	0–100	$2.7 \pm 0.5$	22.45	–	4.76	5.23	5.27
$x_0$ (m)	0–5	$0.2 \pm 0.0$	0.21	0.19	0.2	–	0.18
$h$ (m)	0–60	$4.6 \pm 0.0$	4.69	4.69	4.7	4.84	4.82
$w$ (m)	0–200	$1.8 \pm 0.3$	–	–	–	–	–
Error		$1.5 \times 10^{-5}$	–	–	$4.6 \times 10^{-5}$	–	–

**Acknowledgement** I would like to thank the Editor-in-Chief, Associate Editor and two anonymous reviewers for their comments which have helped to improve the work. This work is a result of a modeling approach in connection with the prospective proposal on the interpretation of mineral exploration study for submission to the Institute of Eminence (IoE) research grant, BHU.

### Compliance with ethical standards

**Conflict of interest** On behalf of all authors, the corresponding author states that there is no conflict of interest.

### References

- Abdelrahman EM (1990) Discussion on “A least-squares approach to depth determination from gravity data” by Gupta O. P. *Geophysics* 55(3):376–378
- Abdelrahman EM, El-Araby TM (1993) A least-squares minimization approach to depth determination from moving average residual gravity anomalies. *Geophysics* 59:1779–1784
- Abdelrahman EM, Sharafeldin SM (1995) A least-squares minimization approach to shape determination from gravity data. *Geophysics* 60:589–590
- Abdelrahman EM, Essa KS (2013) A new approach to semi-infinite thin slab depth determination from second moving average residual gravity anomalies. *Explor Geophys* 44:185–191
- Abdelrahman EM, Essa KS (2015) Three least-squares Minimization Approaches to Interpret Gravity Data Due to Dipping Faults. *Pure Appl Geophys* 172:427–438
- Abdelrahman EM, El-Araby TM, Essa KS (2003) Shape and depth solutions from third moving average residual gravity anomalies using the window curves method: Kuwait. *J Science Eng* 30:95–108
- Abdelrahman EM, Bayoumi AI, El-Araby HM (1991) A least-squares minimization approach to invert gravity data. *Geophysics* 56:115–118
- Abdelrahman EM, Gobashy MM, Essa K, Abo-Ezz ER, El-Araby TM (2016) A least-squares minimization approach to interpret gravity data due to 2D horizontal thin sheet of finite width. *Arab J Geosci* 9:515
- Abedi M, Afshar A, Ardestani VE, Norouzi GH, Lucas C (2010) Application of various methods for 2D inverse modeling of residual gravity anomalies. *Acta Geophys* 58(2):317–336
- Al-Garni MA (2018) Mathematical analysis of gravity anomalies due to an infinite sheet-like structure. *Arab J Geosci* 11:138
- Anderson NL, Essa KS, Elhussein M (2020) A comparison study using particle swarm optimization inversion algorithm for gravity anomaly interpretation due to a 2D vertical fault structure. *J Appl Geophys*. <https://doi.org/10.1016/j.jappgeo.2020.104120>
- Asfahani J, Tlas M (2012) Fair function minimization for direct interpretation of residual gravity anomaly profiles due to spheres and cylinders. *Pure Appl Geophys* 169:157–165
- Atchuta Rao D, Ram Babu HV, Venkata Raju DC (1985) Inversion of gravity and magnetic anomalies over some bodies of simple geometric shape. *Pure Appl Geophys* 123:239–249
- Balkaya C, Ekinci YL, Gokturkler G, Turan S (2017) 3D non-linear inversion of magnetic anomalies caused by prismatic bodies using differential evolution algorithm. *J Appl Geophys* 136:372–386
- Biswas A (2015) Interpretation of residual gravity anomaly caused by a simple shaped body using very fast simulated annealing global optimization. *Geosci Front* 6(6):875–893
- Biswas A (2016) Interpretation of gravity and magnetic anomaly over thin sheet-type structure using very fast simulated annealing global optimization technique. *Model Earth Syst Environ* 2(1):30
- Biswas A, Parija MP, Kumar S (2017) Global nonlinear optimization for the interpretation of source parameters from total gradient of gravity and magnetic anomalies caused by thin dyke. *Ann Geophys* 60(2):1–17
- Blakely RJ, Simpson RW (1986) Approximating edges of source bodies from magnetic or gravity anomalies. *Geophysics* 51:1494–1498
- Ekinci YL, Balkaya C, Gokturkler G, Turan S (2016) Model parameter estimations from residual gravity anomalies due to simple-shaped sources using differential evolution algorithm. *J Appl Geophys* 129:133–147
- Ekinci YL, Balkaya Ç, Göktürkler G (2019) Parameter estimations from gravity and magnetic anomalies due to deep-seated faults: differential evolution versus particle swarm optimization. *Turk J Earth Sci* 28:860–881
- Ekinci YL, Balkaya Ç, Göktürkler G (2020) Global optimization of near-surface potential field anomalies through metaheuristics. In: Biswas A, Sharma SP (eds) *Advances in modeling and interpretation in near surface geophysics*. Springer, Cham, pp 155–188
- Elawadi E, Salem A, Ushijima K (2001) Detection of cavities from gravity data using a neural network. *Explor Geophys* 32:75–79
- Essa KS (2014) New fast least-squares algorithm for estimating the best-fitting parameters of some geometric-structures to measured gravity anomalies. *J Adv Res* 5:57–65
- Essa KS, Munschy M (2019) Gravity data interpretation using the particle swarm optimization method with application to mineral exploration. *J Earth Syst Sci* 128:123
- Essa KS, Géraud Y (2020) Parameters estimation from the gravity anomaly caused by the two-dimensional horizontal thin sheet applying the global particle swarm algorithm. *J Petroleum Sci Eng* 193:107421
- Fernández-Martínez JL, Fernández-Muñiz Z, Cerenea A, Pallero LG, DeAndrés-Galiana EJ, Pedruelo-González LM, Álvarez O, Fernández-Ovies FJ (2020) How to deal with uncertainty in inverse and classification problems. In: Biswas A, Sharma S (eds) *Advances in modeling and interpretation in near surface geophysics*. Springer, Cham, pp 401–414

- Geldart LP, Gill DE, Sharma B (1966) Gravity anomalies of two-dimensional faults. *Geophysics* 31:372–397
- Grant FS, West GF (1965) Interpretation theory in applied geophysics. McGraw Hill Book Co, New York
- Green R (1976) Accurate determination of the dip angle of a geological contact using the gravity method. *Geophys Prospect* 24:265–272
- Gupta OP (1983) A least-squares approach to depth determination from gravity data. *Geophysics* 48:360–375
- Hartmann RR, Teskey D, Friedberg I (1971) A system for rapid digital aeromagnetic interpretation. *Geophysics* 36:891–918
- Hinze WJ, von Frese RRB, Saad AH (2013) Gravity and magnetic exploration: principles, practices and applications. Cambridge University Press, New York
- Holstein H, Fitzgerald D, Anastasiades C (2010) Gravi-magnetic anomalies of uniform thin polygonal sheets. In: 72nd European association of geoscientists and engineers conference and exhibition incorporating SPE EUROPEC, Barcelona, Spain, pp 14–17
- Ingber L, Rosen B (1992) Genetic algorithms and very fast simulated reannealing: a comparison. *Math Comput Model* 16(11):87–100
- Jain S (1976) An automatic method of direct interpretation of magnetic profiles. *Geophysics* 41:531–541
- Kara I, Hoskan N (2016) An easy method for interpretation of gravity anomalies due to vertical finite lines. *Acta Geophys* 64:2232–2243
- Khalil MA, Santos FM, Farzamian M (2014) 3D gravity inversion and Euler deconvolution to delineate the hydro-tectonic regime in El-Arish area, northern Sinai Peninsula. *J Appl Geophys* 103:104–113
- Khalil MA, Santos FM, Farzamian M, El-Kenawy A (2015) 2-D Fourier transform analysis of the gravitational field of Northern Sinai Peninsula. *J Appl Geophys* 115:1–10
- Klingele EE, Marson I, Kahlem HG (1991) Automatic interpretation of gravity gradimetric data in two dimensions: vertical gradient. *Geophys Prospect* 39:407–434
- Lines LR, Treitel S (1984) A review of least-squares inversion and its application to geophysical problems. *Geophys Prospect* 32:159–186
- Mackleod IN, Jones K, Dai TF (1993) 3-D analytical signal in the interpretation of total magnetic field data at low magnetic latitudes. *Explor Geophys* 24:679–688
- Marson I, Klingele EE (1993) Advantages of using the vertical gradient of gravity for 3-D interpretation. *Geophysics* 58:1588–1595
- Mehanee SA (2014) Accurate and efficient regularized inversion approach for the interpretation of isolated gravity anomalies. *Pure Appl Geophys* 171:1897–1937
- Mehanee SA, Essa KS (2015) 2.5D regularized inversion for the interpretation of residual gravity data by a dipping thin sheet: numerical examples and case studies with an insight on sensitivity and non-uniqueness. *Earth Planets Space* 67:130
- Mohan NL, Anandababu L, Roa S (1986) Gravity interpretation using the Mellin transform. *Geophysics* 51:114–122
- Mosegaard K, Tarantola A (1995) Monte Carlo sampling of solutions to inverse problems. *J Geophys Res* 100(B7):12431–12447
- Nabighian MN (1972) The analytical signal of two-dimensional magnetic bodies with polygonal cross-section: its properties and use for automated anomaly interpretation. *Geophysics* 37:507–517
- Nabighian MN (1974) Additional comments on the analytical signal of two-dimensional magnetic bodies with polygonal cross-section. *Geophysics* 39:85–92
- Nettleton LL (1942) Gravity and magnetic calculation. *Geophysics* 7:293–310
- Nettleton LL (1976) Gravity and magnetic in oil prospecting. McGraw Hill Book Co, New York
- Odegard ME, Berg JW (1965) Gravity interpretation using the Fourier integral. *Geophysics* 30:424–438
- Pick M, Picha J, Vyskocil V (1973) Theory of the earth's gravity field. Academia Publishing House of the Czechoslovak Academy of Sciences, Prague, p 538
- Reid AB, Allsop JM, Granser H, Millett AJ, Somerton IW (1990) Magnetic interpretation in the three dimensions using Euler deconvolution. *Geophysics* 55:80–91
- Roy A (1966) The method of continuation in mining geophysical interpretation. *Geoexploration* 4:65–83
- Sen MK, Stoffa PL (1996) Bayesian inference, Gibbs sampler and uncertainty estimation in geophysical inversion. *Geophys Prospect* 44:313–350
- Sen MK, Stoffa PL (2013) Global optimization methods in geophysical inversion. Cambridge University Press, Cambridge
- Sharma SP (2012) VFSARES—a very fast simulated annealing FORTRAN program for interpretation of 1-D DC resistivity sounding data from various electrode array. *Comput Geosci* 42:177–188
- Sharma SP, Biswas A (2013) Interpretation of self-potential anomaly over 2D inclined structure using very fast simulated annealing global optimization—an insight about ambiguity. *Geophysics* 78(3):WB3–WB15
- Sharma B, Geldart LP (1968) Analysis of gravity anomalies of two-dimensional faults using Fourier transforms. *Geophys Prospect* 16:77–93
- Shaw RK, Agarwal BNP (1990) The application of Walsh transforms to interpret gravity anomalies due to some simple geometrically shaped causative sources: a feasibility study. *Geophysics* 55:843–850
- Siegel HO, Winkler HA, Boniwel JB (1957) Discovery of the Mobern Copper Ltd. sulphide deposit, Noranda Mining District, Quebec. In: Methods and case histories in mining geophysics. Commonwealth Mining Metallurgical Congress, 6th, Vancouver, 1957, pp 237–245
- Singh A, Biswas A (2016) Application of global particle swarm optimization for inversion of residual gravity anomalies over geological bodies with idealized geometries. *Nat Resour Res* 25(3):297–314
- Skeels DC (1963) An approximation solution of the problem of maximum depth in gravity interpretation. *Geophysics* 28:724–735
- Srivastava S, Agarwal BNP (2010) Inversion of the amplitude of the two-dimensional analytic signal of the magnetic anomaly by the particle swarm optimization technique. *Geophys J Int* 182:652–662
- Srivastava S, Datta D, Agarwal BNP, Mehta S (2014) Applications of Ant Colony Optimization in determination of source parameters from total gradient of potential fields. *Near Surf Geophys* 12:373–389
- Sundararajan N, Srinivas Y, Laxminarayana Rao T (2000) Sundararajan transform—a tool to interpret potential field anomalies. *Explor Geophys* 31:622–628
- Talwani M, Worze JL, Landisman M (1959) Rapid gravity computations for two-dimensional bodies with applications to the Mendocino submarine fracture zone. *J Geophys Res* 64:49–59
- Thompson DT (1982) EULDPH—a new technique for making computer-assisted depth estimates from magnetic data. *Geophysics* 47:31–37
- Tlas M, Asfahani J (2018) Interpretation of gravity anomalies due to simple geometric-shaped structures based on quadratic curve regression. *Contrib Geophys Geod* 48:161–178
- Tlas M, Asfahani J, Karmeh H (2005) A versatile nonlinear inversion to interpret gravity anomaly caused by a simple geometrical structure. *Pure Appl Geophys* 162:2557–2571
- Trivedi S, Kumar P, Parija MP, Biswas A (2020) Global optimization of model parameters from the 2-D analytic signal of gravity and magnetic anomalies. In: Biswas A, Sharma SP (eds) Advances in modeling and interpretation in near surface geophysics. Springer, Cham, pp 189–221



# Koyna earthquakes: a review of the mechanisms of reservoir-triggered seismicity and slip tendency analysis of subsurface faults

Dip Das<sup>1</sup> · Jyotirmoy Mallik<sup>1</sup>

Received: 18 February 2020 / Accepted: 18 June 2020 / Published online: 6 July 2020  
© Institute of Geophysics, Polish Academy of Sciences & Polish Academy of Sciences 2020

## Abstract

Koyna region in western India experienced more than 1,00,000 earthquakes of different magnitudes ( $M \sim 1.0\text{--}6.3$ ) in the past five decades. Earthquakes in this region are believed to be triggered by a change in fluid pressure due to the percolation of the reservoir (Koyna and Warna reservoir) water into the subsurface. A drilling program was set up by the Ministry of Earth Sciences, India and International Continental Scientific Drilling Program (ICDP) to study the deep subsurface lithology, structure, thermal attributes, etc. as the area is covered by  $\sim 950$  m of thick Deccan basalts. This paper reviews all the hypotheses proposed by earlier workers to explain the mechanism of reservoir trigger causing earthquakes and summarizes such theories to a simple generic model. Slip tendency analysis was further carried out based on the proposed model to explain the dependence of fault slip on fault geometry, rock mechanical properties, stress and fluid gradients. Finally, faults at various depths were characterized (favourably oriented, unfavourably oriented and severely misoriented) based on their potential to go into failure.

**Keywords** Koyna seismicity · Warna seismicity · Reservoir-triggered seismicity · Koyna dam · Slip tendency · Depletion constant

## Introduction

Earthquake is one of the most catastrophic natural calamities that pose serious threats to human life, wealth and economy. It is the result of the sudden release of accumulated stresses in the earth's crust. Stress accumulation can occur due to tectonic plate movement or by human activities such as the construction of dams (Koyna, India), increase in pore pressure by fluid injection like the December 8th, 2006 earthquake in Basel, Switzerland (Bachmann et al. 2012; Mignan et al. 2015) and depletion in reservoir pressure during hydrocarbon production as in case of Groningen gas field, Netherlands (van Thienen-Visser and Breunese 2015; Bommer et al. 2017). Stress relaxation occurs either by creating newer faults in the earth's crust or by reactivating the pre-existing fracture surfaces. When frictional strength between pre-existing fracture surfaces is exceeded by the accumulated stress, faulting occurs (Goswami et al. 2017a), which

may lead to an earthquake. Further, the effect of frictional strength and effective normal stresses can be reduced by the incorporation of fluids into the relatively weaker slip planes (Terzaghi 1943; Skempton 1961; Chen and Nur 1992) which is believed to be the primary reason behind the recurrent earthquakes in the Koyna region, of Maharashtra, Western India.

Koyna, situated on the Deccan Volcanic Province (DVP) in the western part of the Indian peninsula, was classified under a "seismically stable zone" till the early sixties (Rao et al. 1969). Recurrent earthquakes in that region began after the impoundment of the Koyna reservoir (Shivajisagar lake) in 1962 (Gupta and Rastogi 1976; Gupta 1992, 2002; Rastogi et al. 1997; Talwani 1997a, b; Chander and Kalpana 1997; Chadha et al. 1997). This type of seismicity is known as "reservoir-triggered seismicity (RTS)" (Gupta et al. 1969, 1997; Guha et al. 1971; Shashidhar et al. 2011; Dixit et al. 2014). Since then, more than 100,000 earthquakes of magnitude greater than 1.0 (Rastogi et al. 1997; Singh and Chadha 2010), about 200 earthquakes of magnitude around 4 (Dixit et al. 2014) and 22 earthquakes of magnitude greater than 5 (Rao and Shashidhar 2016) were recorded and continue

✉ Jyotirmoy Mallik  
jmallik@iiserb.ac.in

<sup>1</sup> Indian Institute of Science Education and Research, Bhopal, MP 462066, India

being recorded. The largest of them was the  $M \sim 6.3$  earthquake on 10th December 1967.

At any particular depth, the overburden stress is supported by the fluid pressure and the horizontal stresses (Rao et al. 1969). Earthquakes can be triggered by increasing the fluid pressure along the pre-existing fracture planes like it happened in Lake Mead, USA (Carder 1945); Hsinfengkiang dam (1962), China; Kariba dam (1963), Zambia-Zimbabwe; Kremasta dam (1966), Greece (Meade 1991; Gupta 2002). The presence of sedimentary strata in the subsurface can further enhance the occurrence of such seismicity. Sedimentary rocks are capable of retaining water in their pore spaces and that can reduce the effective normal stress and enhance the potential of a fault plane to fail. Rao et al. (1969) attributed the negative Bouguer anomaly over the Koyna region to the presence of sedimentary layers under the Deccan volcanic cover but it lacks direct evidence. Since the region is covered by thick Deccan trap, the subsurface lithology, pore pressure variations and geological structures (faults and flexures), etc. are not well understood and documented. Thus, the existing models explaining RTS in the region are mostly based on scientific speculations with little direct “evidence”.

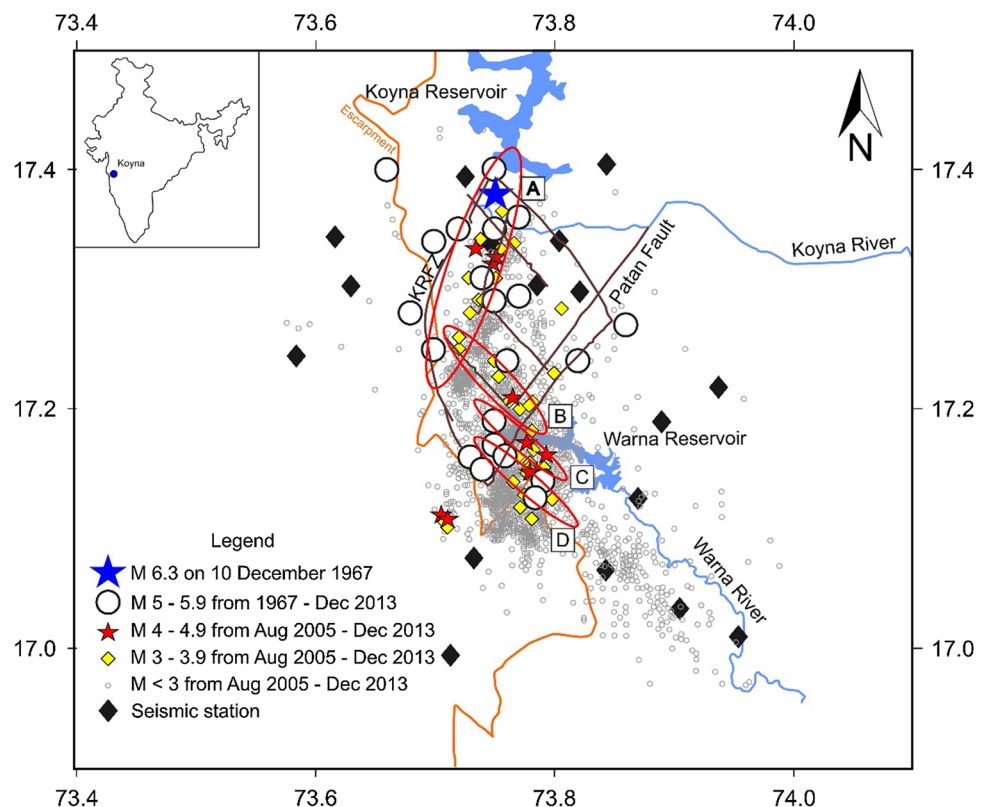
To close such knowledge gaps about the seismotectonics of the region, heat flow structure, geometry and dynamics of the pre-existing weak planes, a scientific deep drilling program was planned in the Koyna area by the Ministry of Earth Sciences, India and Intra Continental Drilling Program

(Gupta et al. 2011, 2015). The drilling revealed the presence of 933 m thick Deccan basalt cover at Koyna underlain by the granitic metamorphic basement with no infra-trappean sedimentary layer (Roy et al. 2013). The cores composed of Deccan volcanics obtained from boreholes show typical signatures of flood basalts and individual pulses that vary in thickness from few metres to tens of metres of flows are identified by the presence of vesicles and amygdules. Well-developed joint sets dissect the massive flood basalt and evidence of fluid movement through them is observed (Roy et al. 2013). Goswami et al. (2017a, b) documented the mechanical properties of underlying granitoids, their porosities and the nature of pre-existing fractures in them. Some of the results from their study are used in the present paper. Below the basalt cover, metamorphosed intermediate/basic rocks are found that correspond to the mid-crustal lithology. The complete absence of granitic crust indicates an erosion of at least  $\sim 15$ – $20$  km of the upper crust before Deccan eruption (Pandey 2016).

### Fault geometry

Rastogi and Talwani (1980) and Langston (1981) suggested the presence of NNE–SSW and NW–SE striking lineaments and Chadha et al. (1997) proposed two NNE–SSW striking faults based on the linear distribution of earthquake epicenters (Fig. 1). The presence of a possible N–S trending shear

**Fig. 1** Distribution of seismicity in the Koyna-Warna seismic zone and associated inferred faults and lineaments. The ellipses represent the spatial clustering of seismic activity along the region. Koyna River Fault Zone (KRFZ) and Patan fault mark the Western and Eastern boundary of the seismic zone, respectively. Several NW–SE trending lineament intersects the KRFZ and Patan fault. Modified after Gupta et al. (2015) and Yadav et al. (2016)



zone along the Koyna river, buried under Deccan basalt, was also proposed by Kailasam and Murthy (1971). Several NW–SE trending fracture planes are present in the region and these fractures should cut across the proposed NNE–SSW trending faults/lineaments (Talwani 1997b). From fault plane solutions, Tandon and Chaudhury (1968) estimated the orientation of the fault planes to be NNE and dipping  $\sim 66^\circ$  towards the west.

### Seismicity in the Koyna region

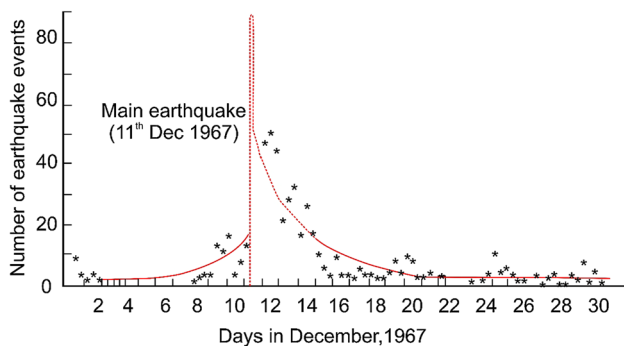
The cluster of earthquake epicenters in the region shifted towards 35 km southward since 1993–1994 after the Warna reservoir dam was built (Rastogi et al. 1997). Depending on the distribution of earthquake epicenters, four well-defined clusters of epicenters were proposed and in recent times further enhancement in seismicity is observed near clusters B, C and D (Fig. 1) (Gupta et al. 2011). In cluster-A, the maximum focal depth of earthquakes are up to about 10 km (Rai et al. 1999; Gupta et al. 2015) and in clusters B, C and D it is up to 8 km (Rai et al. 1999; Dixit et al. 2014 and Gupta et al. 2015). Guha et al. (1968) estimated the foreshock and aftershock pattern of Koyna mainshock (Fig. 2). This type of pattern is similar to Mogi's Type-II model (Gupta et al. 1969) that explains the occurrence of the smaller number of earthquakes before and after the mainshock due to the heterogeneous structure and stress distribution in different lithological units (Mogi 1963). Tandon and Chaudhury (1968) hypothesized a sinistral strike-slip fault responsible for the Koyna mainshock. Langston (1976), on the contrary, envisaged the presence of a sinistral oblique-slip fault dipping towards east.

Rastogi and Talwani (1980) suggested a sinistral strike-slip fault along NNE, near the Koyna reservoir and normal faulting along NW direction about 20 km southward from Koyna reservoir. Mandal et al. (1998) studied the stress drop related to the Koyna earthquakes and suggested two depth

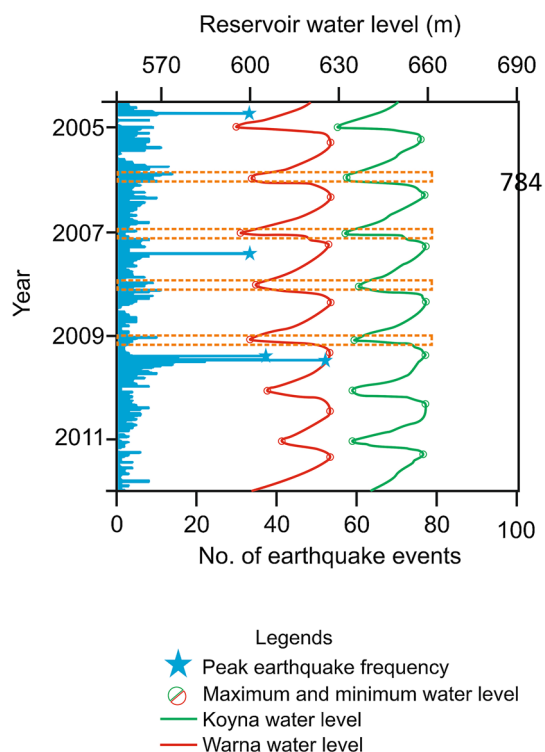
zones, (1) 0–1 km and (2) 5–13 km which show the release of high seismic energy and large stress drops. They argued that the larger stress drop at shallower depth is due to incremental stress in sub-hydrostatic conditions caused by the development of Koyna and Warna reservoirs. At the deeper horizons, this large stress drop is caused by the incremental stress of super-hydrostatic pore fluid pressure diffusion from the reservoirs. Sarma et al. (2004) conducted magnetotelluric surveys on the Koyna region and proposed a 1D model with a moderately conductive vertical layer extended up to a depth of 8 km (seismogenic depth). This low resistive layer could be attributed to the Koyna fault, responsible for the majority of the earthquakes.

### Water-level fluctuations in the Koyna and Warna reservoir and related seismicity

The Koyna region experiences 5000 mm annual rainfall (Yadav et al. 2016). Guha et al. (1966, 1968) correlated the increase in seismic activities to the increase in water level in the reservoir with a time lag. Gupta et al. (1972) further strengthened the argument of relating the occurrence of an earthquake with the increase in water level in the dams, duration of loading, maximum water level and the duration of maximum water level. Gupta (1983) suggested that loading is not the only reason for the earthquakes  $M \geq 5$ . Talwani et al. (1996) found enhanced seismic activity in the region when the highest water level in the reservoir exceeds the previous year maxima. Gupta (2001) studied the correlation between the change in water level in the reservoir and seismicity. He found two seismologically active phases during a year, one is during the pre-monsoon (September) and another is during post-monsoon (February). Pandey and Chadha (2003) studied the relation between the monthly water-level fluctuations and mean monthly strain factors (Lee and Wolf 1998) in Koyna and Warna reservoirs for earthquakes  $M \geq 3$ . This study reveals periodicity in seismic energy release correlated to the annual filling and draining of the reservoirs with a delay of 1 month, until 1996. This delay is related to the diffusion of reservoir water through the pre-existing faults and fractures which increases the pore fluid pressure of the critically stressed medium, facilitating earthquakes. Yadav et al. (2015) also performed a similar study where they compared the periodicity of seismic events with water-level changes in the region. They showed that the number of seismic events increases with the water-level rising in both Koyna and Warna reservoir (Fig. 3). One important thing is to be noticed is that when the water level is at the lowest level, the number of seismic events comparatively increases (Fig. 3, Telesca 2010; Yadav et al. 2015). The majority of earthquake originates from 1 to 7 km depth (Goswami et al. 2017a) and a good correlation was found between the well-level fluctuations and earthquake (Chadha



**Fig. 2** Distribution of foreshock and aftershock in December for Koyna main earthquake (December 10, 1967). Modified after Guha et al. (1968) and Yadav et al. (2016)



**Fig. 3** Water-level fluctuation in Koyna and Warna reservoir from January 2005 to June 2012 and associated earthquake frequency. The dotted boxes mark the lowest levels of water and associated earthquake frequency for the same period. Modified after Yadav et al. (2015)

et al. 1997; Grecksch et al. 1999; Gupta et al. 1999). Rao and Shashidhar (2016) have prepared a catalog of 50 focal mechanism solutions for earthquakes  $M \geq 3.6$ .

In this review, we have summarized and compared the existing models used to explain the earthquakes in Koyna region and tried to find out whether these earthquakes have natural tectonic origin or reservoir-triggered mechanism. We have then summarised the possible models to a generic hypothesis that is based on simple Mohr–Coulomb failure criteria. Using existing data available in the literature on fault attributes, fluid gradient, stress gradient, material properties, we then tried to delineate the “favourably oriented” faults, “unfavourably oriented” faults and “severely misoriented” faults (Sibson 1985) for each fault category. We also discussed “reactivation-tendency factors” or “slip tendency” (Morris et al. 1996; Lisle and Srivastava 2004; Tong and Yin 2011) to quantify the potential of a particular type of fault to acquire slip. We finally tried to shed light on some aspects of periodic seismicity with the loading and unloading cycle of the reservoirs and discussed the limitations of the proposed generic model. We believe this model can also give some theoretical explanations about the future seismic activities.

The next section summarises the existing models proposed by earlier researchers to explain the reservoir-triggered seismicity in the Koyna region.

### Existing models of reservoir-induced seismicity in Koyna region

Talwani (1995) proposed an “intersection model” that suggests the intersection between two or more NW–SE trending faults and major NNE–SSW/N–S features. Such intersections could be responsible for the stress build-up. Stress perturbation could have occurred due to the pore pressure changes caused by the annual loading in the Koyna and Warna reservoirs. He also suggested that the changes in pore fluid pressure could take place due to diffusion with a time lag of 6–8 weeks. This time delay is consistent with the measured rock permeability. Rajendran and Harish (2000) suggested that the Koyna fault is mature and gets weakened by the annual loading cycle of the Koyna reservoir. The failure occurs in this fault due to small changes in stress under high fluid pressure. Talwani (2000) proposed that the rocks in the region are critically stressed and a small change in strength can cause a large earthquake. Pandey and Chadha (2003) simulated the 2D diffusion of pore fluid pressure to study the “fluid pressure diffusion and its relationship with Koyna seismicity”. They assumed a vertical fault and found that it facilitates diffusion. The excess pore pressure can reach up to a depth of 6–8 km. The depth distribution of Koyna earthquakes is consistent with this model. Durá-Gómez and Talwani (2010) validated this model and suggested that fluid pressure above a threshold value causes slippage along the fault and fracture planes. They also suggested that the excess fluid pressure at hypocentral depth along the saturated NE–SW and NW–SE striking faults and fractures is due to the water-level changes in the reservoir. Some researchers, although, (Gahalaut et al. 2004; Srinagesh and Rajagopala Sarma 2005) related the seismicity in Koyna with tectonism and not with the reservoir trigger mechanism.

Significant research was carried out to understand the primary causes responsible for the earthquakes in the region, and robust attempts were made to provide a comprehensive earthquake model for Koyna. Despite that, a proper model that explains all the parameters associated with the earthquakes is lacking. The reason could be the absence of proper crustal velocity model (Yadav et al. 2016), contradictory proposals on fault geometry, lack of understanding about the in situ pore pressure values up to the hypocentral depth, etc.

In summary, scientists believe that the Earthquakes in the Koyna region are mostly due to tectonic stresses but are triggered or aided by the additional stresses provided by the reservoir-triggered pore pressure changes (Gahalaut et al. 2004; Srinagesh and Rajagopala Sarma 2005; Sarkar and

Sain 2017). However, the exact mechanism of such a trigger is not properly understood or documented. The common belief of enhanced pore pressure decreasing the effective stresses leading to differential stress envelope intersecting the failure criteria and causing major brittle failure does not hold good here due to the lack of porous sedimentary formations as proved by deep drilling (Roy et al. 2013; Gupta et al. 2015). Therefore, we propose a simple hypothesis that can explain the mechanism of such a trigger.

## Generic model of Koyna reservoir-triggered seismicity

We propose that the normal stresses which hold fault blocks together (other than friction and asperities between the fault blocks) decrease due to the fluid pressure increase in the fault zone leading to an increase in “slip tendency” which eventually reactivates an already critically stressed fault. The increase in pore pressure happens due to the percolation of surface/rainwater through the fracture system associated with the faults. The overlying Koyna reservoir provides excess fluid pressure, thus facilitating the occurrence of earthquakes.

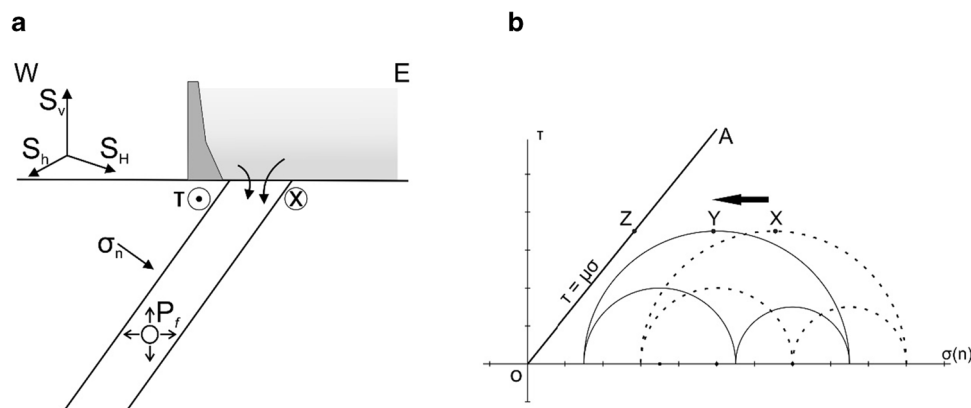
Now, let us consider an inclined, possibly highly dipping, as mentioned earlier (Tandon and Chaudhury 1968), fault zone (Fig. 4a) where the secondary permeability of the fault rock was enhanced by intense fracturing. Let us also assume that the faults already have a higher (> 80–90%) “slip tendency” than a completely inactive or dead fault in a passive

tectonic region (< 30%). The argument in favour of such assumptions is as follows:

In the dilatational step over zone between strike-slip faults, the normal stress decreases and the shear stress increases (Gahalaut and Gahalaut 2008). This favours in normal fault movement. Moreover, Gahalaut et al. (2004) showed that the NNE–SSW left-lateral strike-slip faults are reactivated by N–S compression of the Indian plate and NNW–SSE normal faults are further reactivated by the motion of such strike-slip faults through stress transfer. Thus, the major fault orientations remain “sub-critically” or “critically” stressed either due to the regional stress or the influence of neighbouring faults.

Thus, we think that both the assumptions are quite valid for Koyna based on the findings of previous workers (Tandon and Chaudhury 1968; Gahalaut et al. 2004; Srinagesh and Rajagopala Sarma 2005; Gahalaut and Gahalaut 2008; Sarkar and Sain 2017) as discussed above.

Let us also assume that the fault zone is exposed under principal stresses where  $S_v$  denotes vertical,  $S_H$  and  $S_h$  denote the maximum and minimum horizontal stresses, respectively. Sharma and Mall (1998) have calculated the lithostatic gradient, hydrostatic gradient as well as the fluid pressure gradient for the Koyna region. They proposed a strike-slip stress regime where  $S_H > S_v > S_h$  and this is also a valid assumption as predicted from the focal plane solutions of the seismic faults (Rao and Shashidhar 2016). The maximum horizontal compressive stress ( $S_H$ ) direction in that region is N13°W (Gowd et al. 1992; Mandal and Singh 1996; Heidbach et al. 2016). For simplicity, Andersonian homogeneous fault model is considered and according to



**Fig. 4** **a** Hypothetical East–West section of the study area showing the stress directions as well as the effect of the reservoir and subsequent fluid percolation through the left-lateral fault zone. Here in this study,  $S_v$  corresponds to  $\sigma_2$ ,  $S_H$  corresponds to  $\sigma_1$  and  $S_h$  corresponds to  $\sigma_3$ .  $P_f$  is the fluid pressure inside a hypothetical, steeply west-dipping fault and  $\sigma_n$  is the normal stress on the fault plane.  $\tau$  is shear stress on the fault plane and the dot (.) and cross (x) sign indicates the sense of shear along the fault block. The dot (.) indicates that the movement of the block is towards the South (i.e. towards the

observer) and the cross (x) indicates the movement of the block is towards North (i.e. away from the observer). **b** Shows the effect of fluid pressure increase on 3D Mohr's circle. With increasing fluid pressure, the Mohr's circle will shift towards left and the subsequent slip tendency will increase (follow the text for more details). The dotted circle represents the Mohr's circle before increasing fluid pressure and the solid circle represents the condition after increasing fluid pressure



this model the minimum compressive stress ( $S_h$ ) direction would be N257°W. Of course, a resultant of these stresses will act on the fault planes and it will depend on the angle between the fault plane and regional stress tensor ellipsoid.

Let us plot the resolved normal and shear stresses ( $\sigma_n$ ,  $\tau_s$ ) in the Mohr–coulomb space (denoted by point X) where the horizontal axis represents normal stress (compressive stresses in the positive side) and the vertical axis represents shear stress (Fig. 4b). OA line represents the Mohr–coulomb failure criteria ( $\tau = s_0 + \mu\sigma_n$ ) where  $\mu$  is the coefficient of internal friction represented by the slope of line OA. It is to be noted that the basic assumption of this study is that the fault planes are already existing in nature and the cohesion ( $s_0$ ) is zero in an already fractured plane.

Now, let us assume that after heavy rainfall the reservoirs (Koyna and Warna) above are full of water and creates enough water head so that water can easily percolate through the soil to the fracture and ultimately to the fault zone. This will increase the fluid pressure or pore pressure of the fault rock where porosity and permeability are mostly attributed by the connected fracture system. Due to an increase in pore pressure, which will act opposite to the normal stresses, effective normal stresses will decrease. Please note that the total shear stress on the fault plane is not changed. As a result of decreased effective normal stress, point X will move to a new position towards its left along the horizontal axis, say to point Y. Now as the failure envelope has a positive slope, the slip tendency will increase. With the further addition of pore pressure point, Y may move to point Z and eventually lead to failure.

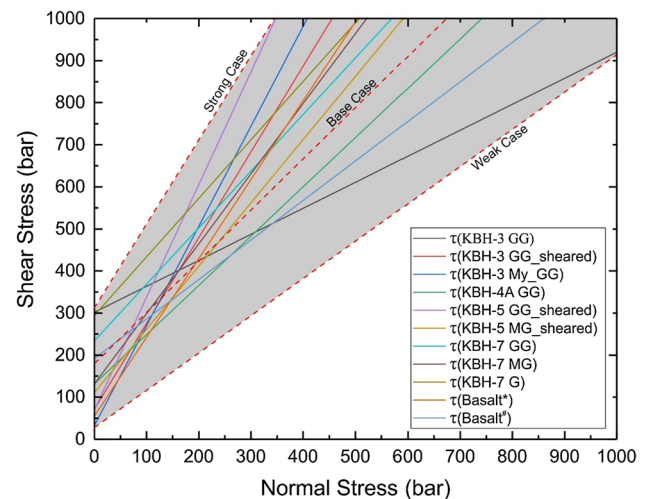
Based on the above model, this study enables us to find out the faults that are prone to slip by the increase in stress or by reduction of effective normal stress on the fault plane due to fluid pressure (Figs. 7, 8). In this study, we have excluded the depths which are beyond the observed hypocentral depth, i.e. 10 km (Rai et al. 1999; Gupta et al. 2015).

The vertical stress ( $S_v$ ) gradient i.e. the lithostatic stress gradient for the region is 0.025 MPa/m (Sharma and Mall 1998). Gowd et al. (1996) calculated the gradient of maximum ( $S_H$ ) and minimum ( $S_h$ ) horizontal stresses from world hydrofrac data and they suggested values like 0.029 MPa/m and 0.0135 MPa/m, respectively.

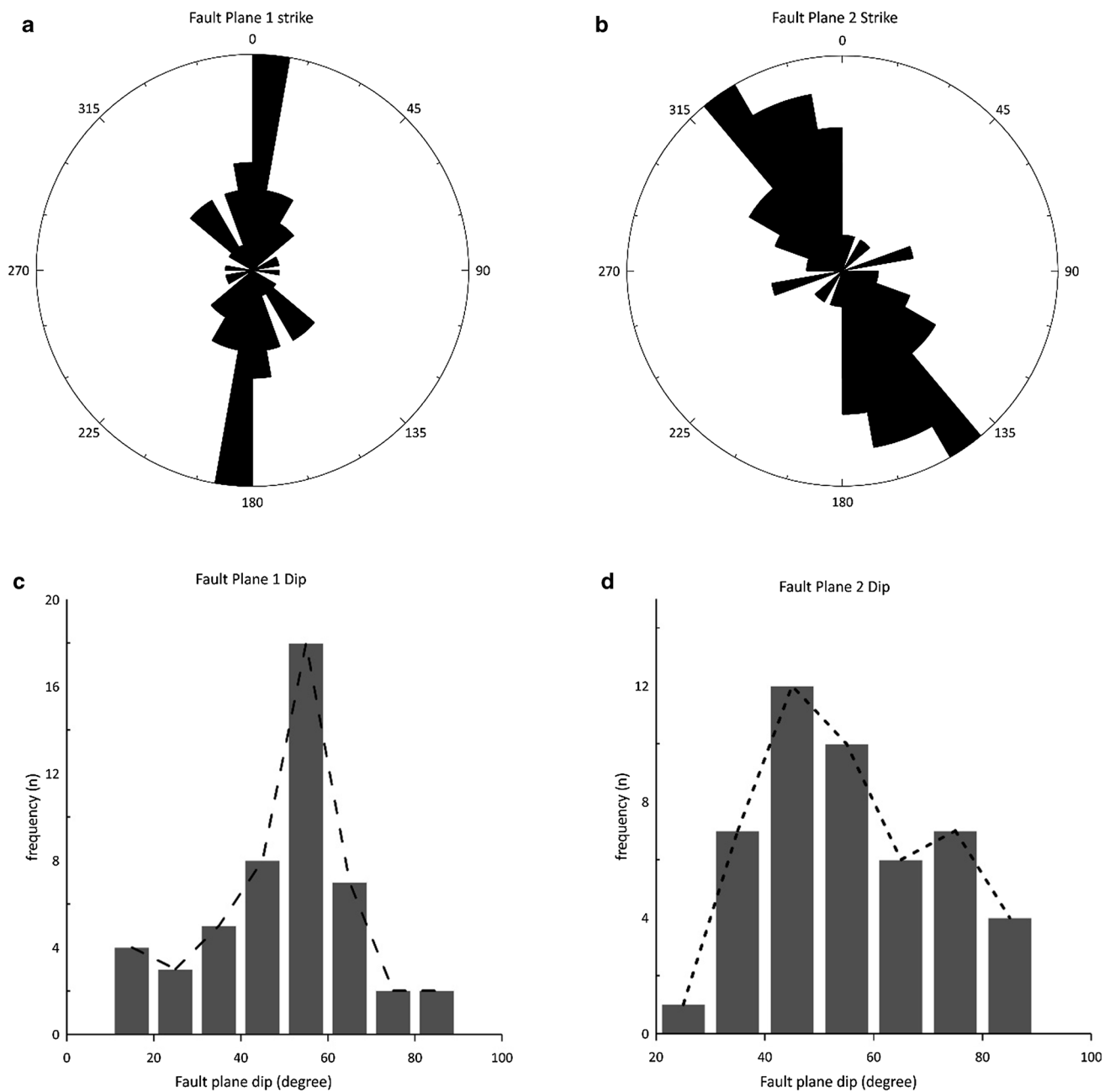
The fluid pressure gradient is considered to be a normal hydrostatic gradient of freshwater which is 0.00979 MPa/m for the Koyna region (Sharma and Mall 1998). The basic assumption for slip tendency calculation is that the pre-existing fault plane has a zero cohesion and the failure envelope is dependent only upon the coefficient of internal friction ( $\mu$ ). Byerlee (1978), based on experimental data, has proposed that  $\mu$  is independent of rock type and a value of 0.85 can be used for normal stress up to 200 MPa. The value of the coefficient of internal friction is highly variable in earth's crust and to understand the effect of  $\mu$  on slip tendency, a

total four different values of  $\mu$  has been chosen for this study. Apart from the value proposed by Byerlee (1978), three other values of  $\mu$  are taken from experimental data by Goswami et al. (2017a) (failure envelopes constrained by the rock strength data obtained from different core samples of KBH-3, KBH-4A, KBH-5 and KBH-7 boreholes). Three hypothetical failure envelopes (Fig. 5) were then constructed [as strong case ( $\varphi = 58.08^\circ$ ,  $\mu = 1.6$ ), base case ( $\varphi = 44.46^\circ$ ,  $\mu = 0.98$ ) and weak case ( $\varphi = 35.53^\circ$ ,  $\mu = 0.71$ )] representing all the aforementioned strength tests.

The orientations of fault planes (Fig. 6) are taken from the focal mechanism solution data published by Rao and Shashidhar (2016). The values of stresses obtained from the above gradients along with the poles of the fault plane attitudes are plotted for different depths on open-source software MohrPlotter (v.3.0) (Allmendinger et al. 2011; Williams et al. 2019; Taghipour et al. 2019). The software returns a Mohr's circle and its corresponding stereonets where the planes are plotted based on its normalized slip tendency values. The warmer colour such as red represents the higher slip tendency, whereas the cooler colour such as blue represents the lesser slip tendency. It is to be noted that a plane can slip when the resolved shear stress on the plane exceeds the frictional resistance of the plane and this frictional resistance is proportional to the normal stress acting



**Fig. 5** This figure shows the failure envelopes measured from core samples collected from the different depth and lithological units of the study area. The samples were obtained from four boreholes named KBH-3, KBH-4A, KBH-5 and KBH-7. The dashed lined represents the three cases chosen for the current analysis; the upper bound envelope as strong case and the lower bound envelope as weak case, respectively. Another in-between case is chosen as the base case. The upper bound and the lower bound lines cover the maximum and minimum extents of the failure envelopes in coordinate space (modified after Goswami et al. 2017a). GG: Granite-gneiss, My\_GG: Mylonitised granite-gneiss; MG: Migmatitic Gneiss, G: Granite, \*: data taken from Burman et al. (2010); #: data are taken from Malik et al. (2017)



**Fig. 6** **a, b** The orientations of the inferred faults from focal mechanism studies. **c, d** The histogram of dip of the inferred faults. Data from Rao and Shashidhar (2016)

along the plane (Jaeger and Cook 1979). Morris et al. (1996) first described the Slip Tendency ( $T_s$ ) to quantify how close a plane is on the verge of failure. It is defined by the ratio of shear stress to normal stress acting on the plane of interest. Slip on a plane is controlled by the cohesive strength ( $S_0$ ) and the coefficient of internal friction ( $\mu$ ) (Morris et al. 1996; Moeck et al. 2009). In this study, we have assumed that the fault planes are already present and they can reactivate by the influence of pore water. Thus, in this case, the cohesion is zero and the slip tendency at failure can be written as:

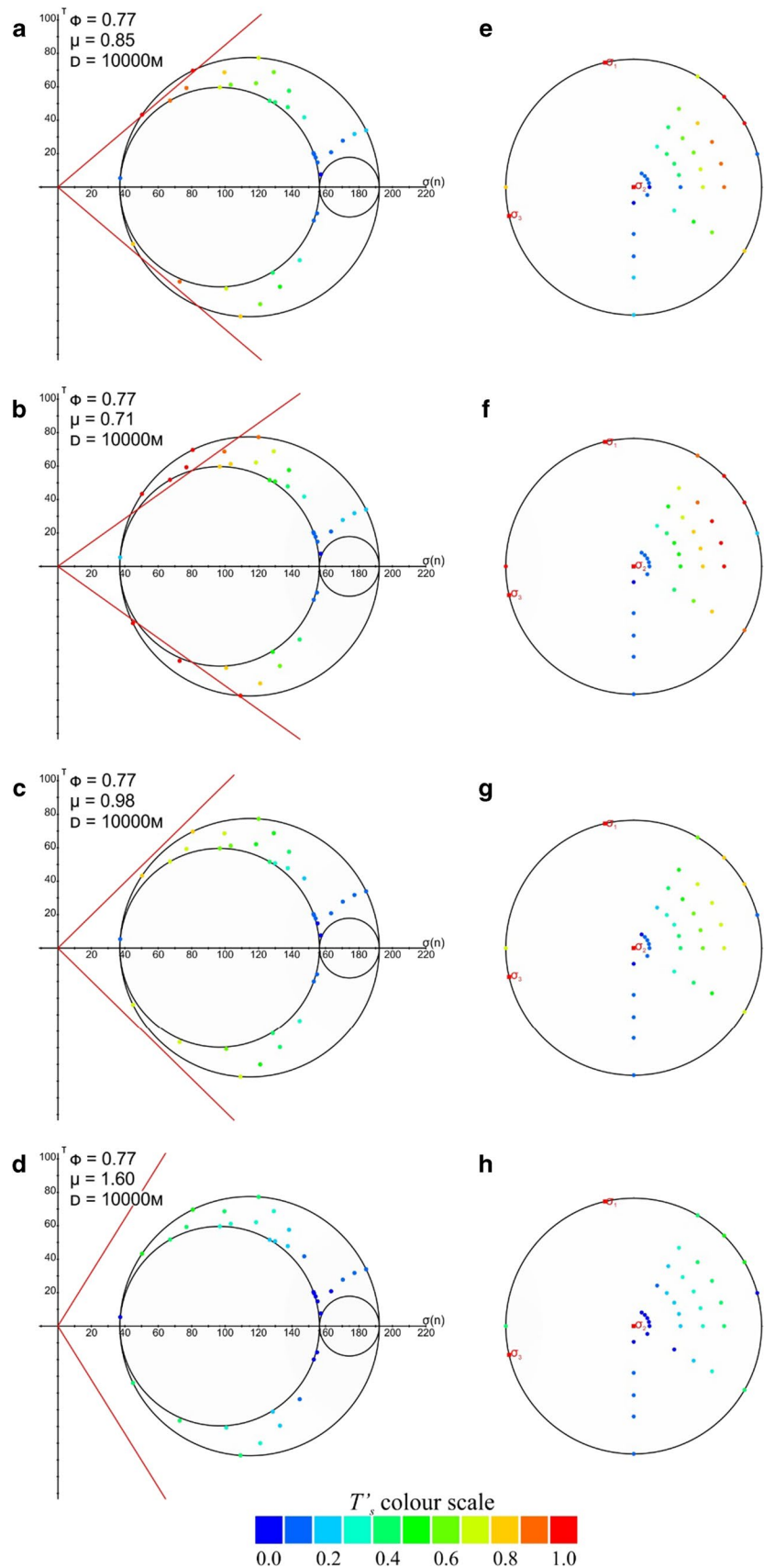
$$T_s = \tau/\sigma \geq \mu \quad (1)$$

Equation 1 can be rewritten as normalized slip tendency index (Lisle and Srivastava 2004) which varies between 0 and 1:

$$T'_s = T_s/\max(T_s) = \tau/\sigma \tan \varphi = \tau/\mu\sigma \quad (2)$$

where  $\varphi$  is the angle of internal friction and  $\tan\varphi = \mu$ .

**Fig. 7** **a–d** 3D Mohr's circles for the different coefficient of internal frictions at stress conditions similar to 10000 m depth. **e–h** represents the stereonet corresponding to each case. The planes on the Mohr's circle and the poles to the planes on the stereonet are colour indexed according to its normalized slip tendency ( $T'_s$ ). The Mohr's circles and stereonets are plotted in MohrPlotter v 3.0 (Allmendinger et al. 2011)



## Inferences and discussion

### Strike-slip reactivation

Figure 7 represents the Mohr's circles and its corresponding stereonet for each case for a depth of 10,000 m. If  $\mu$  is less (0.71–0.85) the N135°E and N150°E trending vertical faults (Fig. 7a, b, e, f) are highly susceptible to reactivation even in normal hydrostatic gradient. The slip along these faults can be attributed to tectonic causes as no overpressure is required for reactivation. Sharma and Mall (1998) have suggested an overpressure at ~3 km depth and this overpressure is high enough to cause reactivation of almost all the pre-existing faults. This would limit the earthquakes with hypocenters at that depth. Since this is not the case for Koyna as evident from Gupta et al. (2016), we have only used the normal hydrostatic gradient to see the effect of reservoir triggering. For “base case” ( $\mu=0.98$ ), the aforementioned faults show a  $T'_s$  of 0.8 and for “strong case” ( $\mu=1.6$ ) it is 0.5 (Fig. 7c, d, g, h). It is to be noted that this scenario is for 10,000 m depth with a normal hydrostatic gradient. Thus, the faults which are reactivating by strike-slip movement (NNW to NS trending vertical to sub-vertical faults) or are on the verge of reactivation by little overpressure can be attributed to natural tectonic causes. These faults can be said to be “favourably” oriented for reactivation. It is obvious that fluids have a role to play in these reactivations (King et al. 1994; Roeloffs 1996) as the presence of fluid and fluid pressure changes at depths have been reported by many authors (Sharma and Mall 1998; Talwani et al. 1999; Gupta et al. 2000; Kalpana and Chander 2000). As a result, the actual earthquake will happen at a much lesser depth. These faults are already critically stressed. The strike-slip stress regime, the high dip of the reactivating fault planes (Fig. 7) and the above-mentioned point are in agreement with the finding of the previous authors, i.e. these faults are indeed tectonically reactivated (Gahalaut et al. 2004; Srinagesh and Rajagopala Sarma 2005).

In the study area, sinistral strike-slip movement has been reported from NS to NNE trending faults and this is favoured by the plate tectonics. But in this study, the NW–SE trending faults are also showing high potential for strike-slip reactivation. Arora et al. (2018) have proposed a possible connection between Chitradurga Shear Zone (CSZ) and two major NW–SE faults in the study area. They believe that the major NW–SE trending faults are an extension of the CSZ. This can explain our result on the strike-slip movement along the NW–SE trending faults. Since we have chosen hypothetical fault planes and these planes are representative of a particular set, we do not know the exact location of each fault plane. If the bounding faults of the extended CSZ are defined by these NW–SE-oriented faults, the absence of

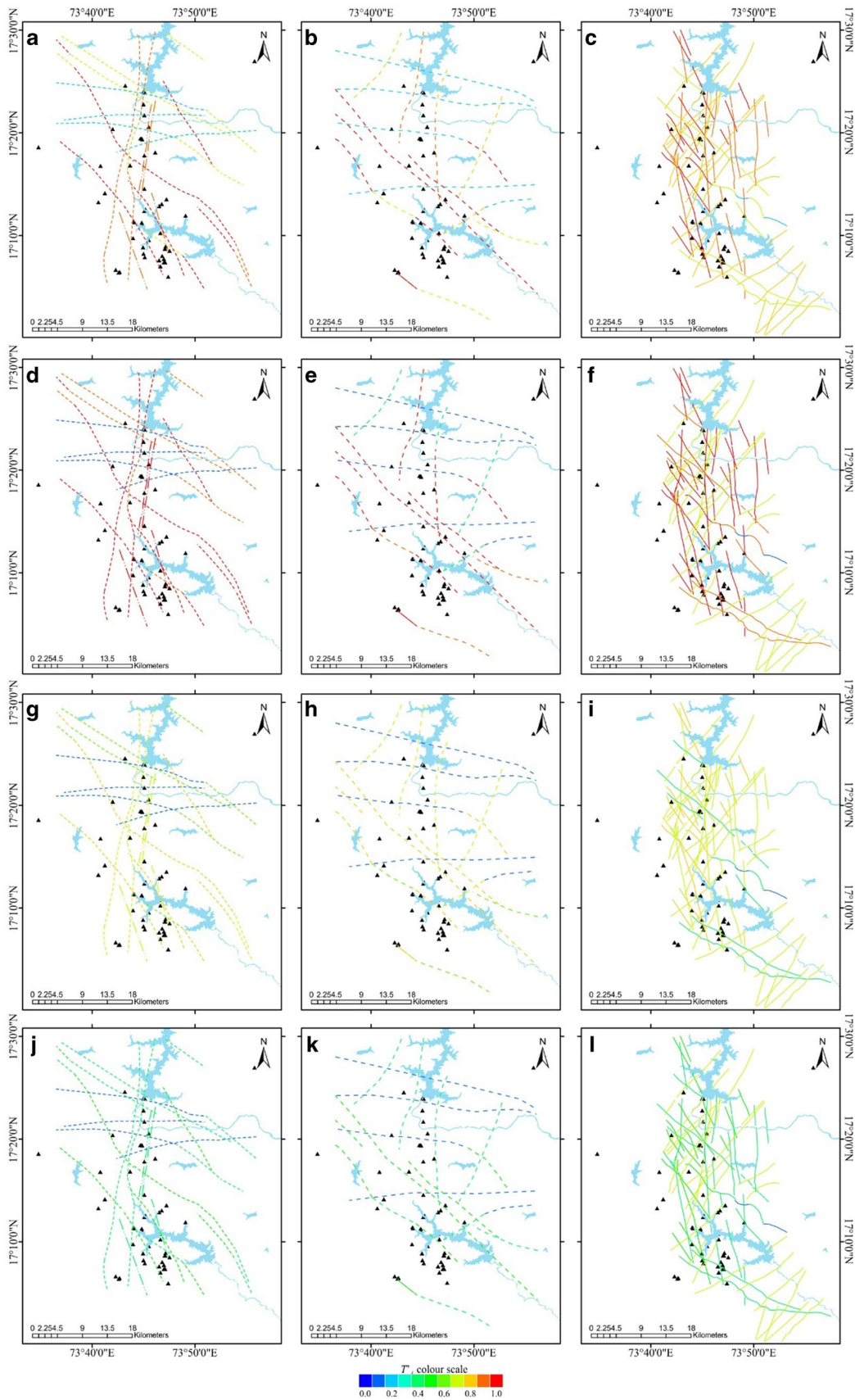
strike-slip earthquake on these faults can be attributed to that fact that in shear zones, seismicity is mainly restricted between the bounding faults and a scattering of epicenters is observed in-between those two faults (Arora et al. 2018). The faults with lower dip are not critically stressed and less likely to be reactivated by tectonic causes rather these would be reactivated if a higher fluid pressure is supplied such as by construction of a reservoir on the surface. The effect of the reservoir impoundment on seismicity is discussed in the latter section of this text.

### Normal reactivation

Aftershock due to normal reactivation of faults is reported by many authors (Gupta 1992; Mandal et al. 2000; Gahalaut and Gahalaut 2008 and references therein). Gahalaut and Gahalaut (2008) have discussed an interesting phenomenon of normal earthquake in compressive settings. They showed that at step over zone related to strike-slip faulting, the faults in the dilatational zone are more likely to be reactivated than the faults in the compressional zone. This is due to the fact that in the dilatational zone the shear stress increases and the normal stress decreases whereas in compressional zone normal and shear stress both increases thus favouring normal reactivation of faults in the dilatational zones. They also showed that static stress increases with the increasing effective coefficient of friction ( $\mu'$ ) in the dilatational zone whereas it decreases with increasing  $\mu'$  in the compressional zone. Thus the lower dip faults in “favourable orientation” are passively stressed by the activity in the strike-slip faults and are likely to be reactivated as normal faults.

### Relationship between fault geometry and reactivation potential

In Fig. 8, the lineation within the upper basaltic unit (Fig. 8a, d, g, j), basement granitoid (Fig. 8b, e, h, k) and on the surface (Fig. 8c, f, i, l) are marked with respect to their  $T'_s$ . Since individual dips for these planes are not known, we have chosen to mark the maximum slip tendency possible for each case. The readers are advised to keep in mind that not all of the faults for a particular scenario and in a particular medium will be reactivated simultaneously. The faults which are showing favourable dip will only be reactivated. Figure 8 only displays all possible faults that can be reactivated and how likely they are to be reactivated. It is evident from Figs. 7 and 8 that the N120°E trending faults are showing moderate  $T'_s$  and these planes can be said “unfavourably” oriented for reactivation. Whereas, the EW-oriented faults are showing very low  $T'_s$  and these planes can be said “severely misoriented” for reactivation and seismicity in these faults is less likely to occur. In order to validate our model, we calculated the slip tendency of existing seismogenic faults



**Fig. 8** Lineament maps of the study area for different coefficient of internal frictions (a–c for  $\mu=0.85$ ; d–f for  $\mu=0.71$ ; g–i for  $\mu=0.98$  and j–l for  $\mu=1.6$ ). a, d, g and j represents lineation in upper Deccan basaltic unit, b, e, h and k represents lineation in the basement granitoids and c, f, i and l represents the surface lineation. The lineaments are coloured based on their maximum  $T'_s$  possible (Fig. 7). Modified after Arora et al. (2017, 2018)

(faults which already slipped and resulted earthquake). The analysis shows that such faults are either “critically” or “sub-critically” stressed depending on the value of frictional coefficient used (Fig. 9). The results corroborates well with the seismogenic faults that slipped during pre- and post-monsoon season due to water-level fluctuations in the Koyna and Warna reservoirs (see the stereonet in Fig. 9). Figure 10a shows the orientation of the “favourable”, “unfavourable” and “severely misoriented” faults with respect to the regional stress field of the area and Fig. 10b the mechanism of failure along the “favourably” oriented faults.

### Effect of reservoir trigger on seismicity

The impoundment of the Koyna reservoir would exert excess pressure on the ground and the percolating fluid will supply excess fluid pressure, favouring the earthquake occurrence. The fracture networks within the lithological units would enhance the permeability (Arora et al. 2017) and facilitate the diffusion of water up to the seismogenic depth. One interesting fact about it is that each earthquake would lead to the creation of fractures within the lithological units which will eventually increase the permeability of the rocks (Gavrilenko et al. 2010). So, the tectonic stress along the probable normal faults are high enough due to strike-slip activity in other faults and a small increase in fluid pressure which corresponds to a lesser depth can trigger an earthquake. Whereas the strike-slip faults are also critically stressed but there is no enhancement of tectonic stresses, thus it will fail at a greater depth compared to the normal faults. This phenomenon can explain the observation by Mandal et al. (1998), i.e. the strike-slip movements are responsible for deeper events (focal depth  $\geq 5$  km) and the normal movements are responsible for shallower (focal depth  $\leq 5$  km) events.

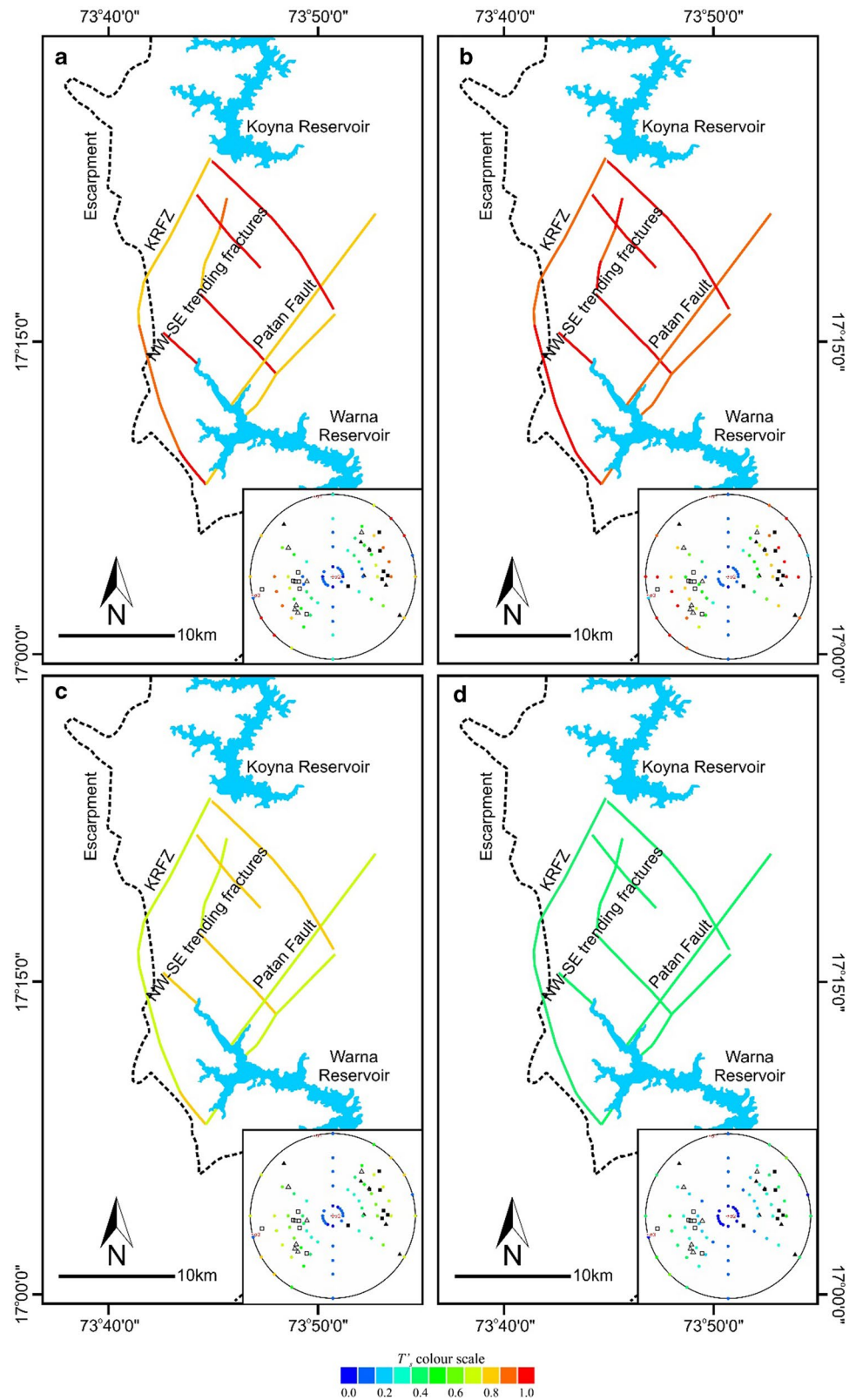
The increase in seismicity during the highest water level in the reservoir can be attributed to the enhanced fluid pressure at depths due to diffusion of water through fracture networks (Chadha et al. 1997; Grecksch et al. 1999; Gupta et al. 1999; Gupta 2001; Pandey and Chadha 2003; Yadav et al. 2015) leading to an increasing  $T'_s$ . However, the little comparative increase in seismic events at the lowest water level can be explained by depletion constant ( $\gamma$ ) (Hillis 2001; Zoback and Zinke 2002; Liu and Harpalani 2014; Atapour and Mortazavi 2018).  $\gamma$  is defined by:

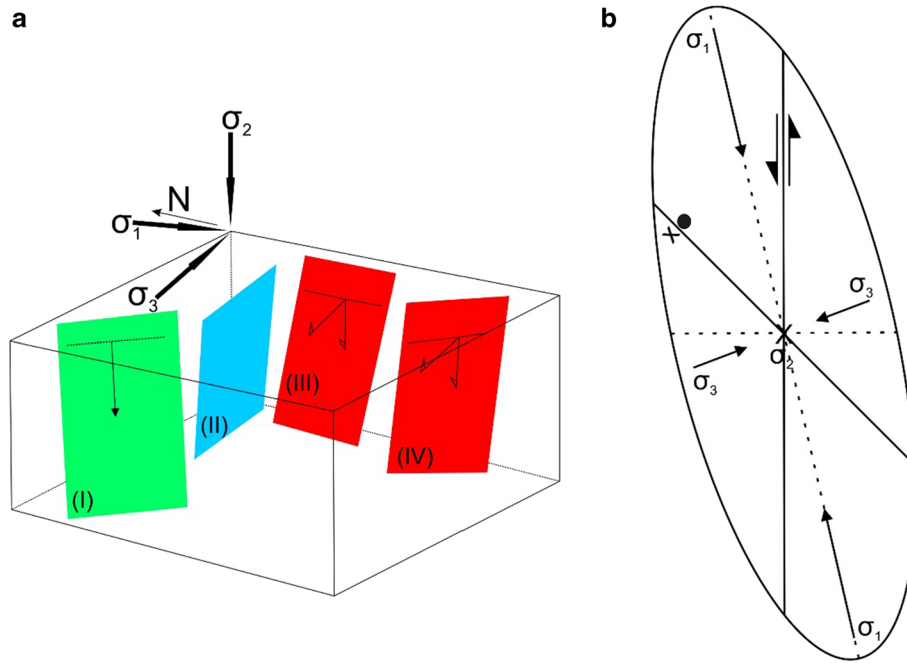
$$\Delta S_h / \Delta P = \gamma \quad (3)$$

where  $\Delta S_h$  is change in minimum horizontal stress and  $\Delta P$  is change in pore pressure. From Eq. 3, it is evident that change in minimum horizontal stress is linearly related to change in pore pressure. So, at the lowest water level, the pore pressure will decrease and it will lead to a decrease in effective minimum horizontal stress. This characteristic of  $S_h$  will increase the radius of Mohr’s circle and the slip tendency of the fault plane will increase. Thus, it will lead to seismic events related to depletion. This depletion-related seismicity is pretty common for the mature oil and gas field (Groningen gas field, The Netherlands), and this is the first time where this phenomenon is used to describe the seismicity related to reservoir triggering.

This model can only predict which faults (based on its strike) are likely to be reactivated but fails to exactly identify the faults geographically until the complete attitude of each fault is available. The specific  $T'_s$  value of an individual fault is overestimated to some extent as representative  $T'_s$  value is assigned to it which could be slightly higher. In this model, westerly dipping faults are chosen (Tandon and Chaudhury 1968; Lee and Raleigh 1969; Sykes 1970; Singh et al. 1975), but this method is insensitive to the dip direction as stress is bi-directional. This model limits itself to identify the exact mechanisms behind the continuous earthquake as well as the exact fluid pressure conditions in the region and the source parameters for earthquakes which originate below 10 km depth. The primary parameters used in this model are fault geometry, stress orientation and magnitude, fluid pressure and rock mechanical properties like cohesion and frictional coefficient. Amongst these parameters, fault geometry is rather well constrained from earthquake focal point solutions. For stress orientation regional  $S_H$  azimuth is used and that is also well constrained from many earthquake focal point solution data. To calculate stress magnitude, we have used standard stress gradient from literature which of course can locally vary depending on depth, fluid pressure and tectonic conditions. We have not considered any overpressure in our model and addition of proper overpressure information can significantly constrain the model and improve outcome. Although, we have used a range in rock mechanical properties from the literature, assignment of specific rock mechanical properties to individual faults is virtually impossible for such a large-scale generic model. Overall, all the modelling parameters are validated from different sources and hence the modelling result should be a close representation of the reality.

**Fig. 9** The slip tendency ( $T'_s$ ) for the known seismogenic faults in the Koyna region for **a**  $\mu=0.85$ , **b**  $\mu=0.71$ , **c**  $\mu=0.98$  and **d**  $\mu=1.6$ . The inset represents stereographic projections (poles) of seismogenic faults and corresponding slip-tendencies (designated by its colour) for corresponding  $\mu$  values. The triangles represent pre-monsoon seismogenic faults and squares represent post-monsoon seismogenic faults (2005–2012, Rao and Shashidhar 2017 and references there in). Solid symbols represent the primary fault planes and the hollow symbols represent the auxiliary fault planes. Modified after Rao et al. (2017)





**Fig. 10** **a** Schematic diagram showing the orientations of “favourably”, “unfavourably” and “severely misoriented” faults for reactivation with respect to the regional stress field. The red coloured planes (III and IV) and the marked dip show the attitude of the plane with favourable orientation, the greener one (I) shows the attitude of unfavourable orientation and the blue (II) one shows the attitude of severely unfavourable fault. Note that the blue (II) one is severely misoriented and it is irrespective of its dip whereas for the other

two dip amount plays a significant role in its reactivation potential (Fig. 7). **b** Shows the regional stress field and the slip mechanism along “favourably” oriented faults. The NE–SW trending faults are likely to be reactivated by strike-slip mechanism and the NW–SE trending faults by right-lateral strike-slip movement but normal movement along these faults are observed (Gahalaut and Gahalaut 2008)

## Conclusions

Our generic model suggests that the faults are already critically stressed due to tectonic forces and the impoundment of the Koyna reservoir further increases the vertical load as well as the fluid pressure. The enhanced fluid pressure at seismogenic depths due to a well-developed fracture network in basement rocks and resulting high diffusivity of the fluids decreases the normal stress on the already critically stressed faults. The NS to NNE trending vertical faults are favourably oriented for reactivation (i.e. showing higher slip tendency) and strike-slip movement can be triggered by little increase in fluid pressure. The strike-slip motion along these faults favours normal faulting along the dilatational step over zones. These normal faulting can be further enhanced by the incorporation of fluids into the fracture system. Henceforth, the earthquakes in this region can be attributed to mainly tectonic causes with the additional effect of the Koyna reservoir. Thus, Koyna stands out to be a perfect example of reservoir-triggered seismicity.

**Acknowledgements** This research work is part of DD’s Doctoral research. DD and JM like to thank Indian Institute of Science Education and Research (IISER) Bhopal for its financial support. The authors

also like to acknowledge Mr. Krishanu Bandhopadhyay for his valuable suggestions. The authors acknowledge the donation of the academic license of MOVE software suite by Petroleum Experts Limited on which the results of the Slip-tendency analysis were verified.

## Compliance with ethical standards

**Conflict of interest** The authors declare that they have no conflict of interest.

## References

- Allmendinger RW, Cardozo N, Fisher DM (2011) Structural geology algorithms: vectors and tensors. Cambridge University Press, Cambridge
- Arora K, Chadha RK, Srinu Y, Selles A, Davuluri S, Smirnov V, Ponomarev A, Mikhailov VO (2017) Lineament fabric from airborne LiDAR and its influence on triggered earthquakes in the Koyna-Warna region, western India. *J Geol Soc India* 90(6):670–677
- Arora K, Srinu Y, Gopinadh D, Chadha RK, Raza H, Mikhailov V, Ponomarev A, Kiseleva E, Smirnov V (2018) Lineaments in Deccan Basalts: the basement connection in the Koyna-Warna RTS region. *Bull Seismol Soc Am* 108(5B):2919–2932
- Atapour H, Mortazavi A (2018) The influence of mean grain size on unconfined compressive strength of weakly consolidated reservoir sandstones. *J Petrol Sci Eng* 171:63–70



- Bachmann CE, Wiemer S, Goertz-Allmann BP, Woessner J (2012) Influence of pore-pressure on the event-size distribution of induced earthquakes. *Geophys Res Lett* 39(9):L09302
- Bommer JJ, Stafford PJ, Edwards B, Dost B, van Dedem E, Rodriguez-Marek A, Kruiver P, van Elk J, Doornhof D, Ntinalexis M (2017) Framework for a ground-motion model for induced seismic hazard and risk analysis in the Groningen gas field, the Netherlands. *Earthq Spectra* 33(2):481–498
- Burman A, Maity D, Sreedeeep S (2010) Iterative analysis of concrete gravity dam-nonlinear foundation interaction. *Int J Eng Sci Technol* 2(4):85–99
- Byerlee J (1978) Friction of rocks. In: *Rock friction and earthquake prediction*. Birkhäuser, Basel, pp 615–626
- Carder DS (1945) Seismic investigations in the Boulder Dam area, 1940–1944, and the influence of reservoir loading on local earthquake activity. *Bull Seismol Soc Am* 35(4):175–192
- Chadha RK, Gupta HK, Kumpel HJ, Mandal P, Rao AN, Kumar N, Radhakrishna I, Rastogi BK, Raju IP, Sarma CSP, Satyamurthy C (1997) Delineation of active faults, nucleation process and pore pressure measurements at Koyna (India). In: *Seismicity associated with mines, reservoirs and fluid injections*. Birkhäuser, Basel, pp 551–562
- Chander R, Kalpana (1997) On categorising induced and natural tectonic earthquakes near new reservoirs. *Eng Geol* 46(2):81–92
- Chen Q, Nur A (1992) Pore fluid pressure effects in anisotropic rocks: mechanisms of induced seismicity and weak faults. *Pure Appl Geophys* 139(3–4):463–479
- Dixit MM, Kumar S, Catchings RD, Suman K, Sarkar D, Sen MK (2014) Seismicity, faulting, and structure of the Koyna-Warna seismic region, Western India from local earthquake tomography and hypocenter locations. *J Geophys Res Solid Earth* 119(8):6372–6398
- Durá-Gómez I, Talwani P (2010) Hydromechanics of the Koyna-Warna region. India. *Pure Appl Geophys* 167(1–2):183–213
- Gahalaut K, Gahalaut VK (2008) Stress triggering of normal aftershocks due to strike slip earthquakes in compressive regime. *J Asian Earth Sci* 33(5–6):379–382
- Gahalaut VK, Kalpana, Singh SK (2004) Fault interaction and earthquake triggering in the Koyna-Warna region, India. *Geophys Res Lett* 31(11)
- Gavrilenko P, Singh C, Chadha RK (2010) Modelling the hydromechanical response in the vicinity of the Koyna reservoir (India): results for the initial filling period. *Geophys J Int* 183(1):461–477
- Goswami D, Akkiraju VV, Misra S, Roy S, Singh SK, Sinha A, Gupta H, Bansal BK, Nayak S (2017a) Rock strength measurements on Archaean basement granitoids recovered from scientific drilling in the active Koyna seismogenic zone, western India. *Tectonophysics* 712:182–192
- Goswami D, Akkiraju VV, Singh SK, Roy S (2017b) Rock strength and elastic properties of basement granitoids from Koyna region, Deccan Volcanic Province, India. *J Geol Soc India* 90(6):783–787
- Gowd TN, Rao SS, Gaur VK (1992) Tectonic stress field in the Indian subcontinent. *J Geophys Res Solid Earth* 97(B8):11879–11888
- Gowd TN, Rao SS, Chary KB (1996) Stress field and seismicity in the Indian shield: effects of the collision between India and Eurasia. In: *Mechanics problems in geodynamics, part II*. Birkhäuser Basel, pp 503–531
- Grecksch G, Roth F, Kumpel HJ (1999) Coseismic well-level changes due to the 1992 Roermond earthquake compared to static deformation of half-space solutions. *Geophys J Int* 138(2):470–478
- Guha SK, Gosavi PD, Padale JG, Marwadi SC (1966) Crustal disturbance in the Shivaji Sagar-Lake area of the Koyna hydroelectric project (Maharashtra, India). In: *3rd Symposium on earthquake engineering*, Univ. Roorkee
- Guha SK, Gosavi PD, Varma MM, Agrawal SP, Padale JG, Marwadi SC (1968) Recent seismic disturbances in the Koyna Hydroelectric Project. Central Water and Power Research Station, Maharashtra
- Guha SK, Gosavi PD, Padale JG, Marwadi SC (1971) An earthquake cluster at Koyna. *Bull Seismol Soc Am* 61(2):297–315
- Gupta HK (1983) Induced seismicity hazard mitigation through water level manipulation at Koyna, India: a suggestion. *Bull Seismol Soc Am* 73(2):679–682
- Gupta HK (1992) *Reservoir induced earthquakes*, vol 64. Elsevier, Amsterdam
- Gupta HK (2001) Short-term earthquake forecasting may be feasible at Koyna, India. *Tectonophysics* 338(3–4):353–357
- Gupta HK (2002) A review of recent studies of triggered earthquakes by artificial water reservoirs with special emphasis on earthquakes in Koyna, India. *Earth Sci Rev* 58(3–4):279–310
- Gupta HK, Rastogi BK (1976) Dams and earthquakes. In: *Developments in geotechnical engineering*, no. 11. Elsevier, Amsterdam
- Gupta H, Narain H, Rastogi BK, Mohan I (1969) A study of the Koyna earthquake of December 10, 1967. *Bull Seismol Soc Am* 59(3):1149–1162
- Gupta HK, Rastogi BK, Narain H (1972) Common features of the reservoir-associated seismic activities. *Bull Seismol Soc Am* 62(2):481–492
- Gupta HK, Rastogi BK, Chadha RK, Mandal P, Sarma CSP (1997) Enhanced reservoir-induced earthquakes in Koyna region, India, during 1993–95. *J Seismol* 1(1):47–53
- Gupta HK, Rao RUM, Srinivasan R, Rao GV, Reddy GK, Dwivedy KK, Banerjee DC, Mohanty R, Satyasaradhi YR (1999) Anatomy of surface rupture zones of two stable continental region earthquakes, 1967 Koyna and 1993 Latur, India. *Geophys Res Lett* 26(13):1985–1988
- Gupta HK, Radhakrishna I, Chadha RK, Kumpel HJ, Grecks G (2000) Pore pressure studies initiated in area of reservoir-induced earthquakes in India. *Eos Trans Am Geophys Union* 81(14):145–151
- Gupta H, Shashidhar D, Mallika K, Rao NP, Srinagesh D, Satyanarayana HVS, Saha S, Naik RTB (2011) Short term earthquake forecasts at Koyna, India. *J Geol Soc India* 77(1):5–11
- Gupta H, Rao NP, Roy S, Arora K, Tiwari VM, Patro PK, Satyanarayana HVS, Shashidhar D, Mallika K, Akkiraju VV, Goswami D (2015) Investigations related to scientific deep drilling to study reservoir-triggered earthquakes at Koyna, India. *Int J Earth Sci* 104(6):1511–1522
- Gupta HK, Arora K, Rao NP, Roy S, Tiwari VM, Patro PK, Satyanarayana HVS, Shashidhar D, Mahato CR, Srinivas KNSSS, Srihari M (2016) Investigations of continued reservoir triggered seismicity at Koyna, India. *Geol Soc Lond Spec Publ* 445(1):151–188
- Heidbach O, Rajabi M, Reiter K, Ziegler M, WSM Team (2016) World stress map database release 2016. GFZ Data Services
- Hillis RR (2001) Coupled changes in pore pressure and stress in oil fields and sedimentary basins. *Pet Geosci* 7(4):419–425
- Jaeger JC, Cook NG (1979) *Fundamentals of rock mechanics*. Chapman and Hall, London, p 593
- Kailasam LN, Murthy BGK (1971) A short note on gravity and seismic investigations in the Koyna area. *Indian J Power River Val Dev Spec Number* 33:27–30
- Kalpana, Chander R (2000) Green's function based stress diffusion solutions in the porous elastic half space for time varying finite reservoir loads. *Phys Earth Planet Inter* 120(1–2):93–101. [https://doi.org/10.1016/S0031-9201\(00\)00146-1](https://doi.org/10.1016/S0031-9201(00)00146-1)
- King GC, Stein RS, Lin J (1994) Static stress changes and the triggering of earthquakes. *Bull Seismol Soc Am* 84(3):935–953
- Langston CA (1976) A body wave inversion of the Koyna, India, earthquake of December 10, 1967, and some implications for body wave focal mechanisms. *J Geophys Res* 81(14):2517–2529
- Langston CA (1981) Source inversion of seismic waveforms: the Koyna, India, earthquakes of 13 September 1967. *Bull Seismol Soc Am* 71(1):1–24

- Lee WHK, Raleigh CB (1969) Fault-plane solution of the Koyna (India) earthquake. *Nature* 223(5202):172–173
- Lee MK, Wolf LW (1998) Analysis of fluid pressure propagation in heterogeneous rocks: implications for hydrologically-induced earthquakes. *Geophys Res Lett* 25(13):2329–2332
- Lisle RJ, Srivastava DC (2004) Test of the frictional reactivation theory for faults and validity of fault-slip analysis. *Geology* 32(7):569–572
- Liu S, Harpalani S (2014) Evaluation of in situ stress changes with gas depletion of coalbed methane reservoirs. *J Geophys Res Solid Earth* 119(8):6263–6276
- Malik A, Chakraborty T, Rao KS, Kumar D, Chandel P, Sharma P (2017) Dynamic response of Deccan Trap basalt under Hopkinson bar test. *Procedia Eng* 173:647–654
- Mandal P, Singh RN (1996) Three-dimensional intraplate stress distributions associated with topography and crustal density inhomogeneities beneath the Deccan Volcanic Province. *Proc Indian Acad Sci Earth Planet Sci* 105(2):143–155
- Mandal P, Rastogi BK, Sarma CSP (1998) Source parameters of Koyna earthquakes, India. *Bull Seismol Soc Am* 88(3):833–842
- Mandal P, Rastogi BK, Gupta HK (2000) Recent Indian earthquakes. *Curr Sci Bangalore* 79(9):1334–1346
- Meade RB (1991) Reservoirs and earthquakes. *Eng Geol* 30(3–4):245–262
- Mignan A, Landtwing D, Kästli P, Mena B, Wiemer S (2015) Induced seismicity risk analysis of the 2006 Basel, Switzerland, Enhanced Geothermal System project: influence of uncertainties on risk mitigation. *Geothermics* 53:133–146
- Moock I, Kwiatak G, Zimmermann G (2009) Slip tendency analysis, fault reactivation potential and induced seismicity in a deep geothermal reservoir. *J Struct Geol* 31(10):1174–1182
- Mogi K (1963) Some discussions on aftershocks, foreshocks and earthquake swarms—the fracture of a semi-infinite body caused by an inner stress origin and its relation to the earthquake phenomena. *Bull Earthq Res Inst* 41:615–658
- Morris A, Ferrill DA, Henderson DB (1996) Slip-tendency analysis and fault reactivation. *Geology* 24(3):275–278
- Pandey OP (2016) Deep scientific drilling results from Koyna and Kilarī earthquake regions reveal why Indian shield lithosphere is unusual, thin and warm. *Geosci Front* 7(5):851–858
- Pandey AP, Chadha RK (2003) Surface loading and triggered earthquakes in the Koyna-Warna region, western India. *Phys Earth Planet Inter* 139(3–4):207–223
- Rai SS, Singh SK, Sarma PR, Srinagesh D, Reddy KNS, Prakasam KS, Satyanarayana Y (1999) What triggers Koyna region earthquakes? Preliminary results from seismic tomography digital array. *Proc Indian Acad Sci Earth Planet Sci* 108(1):1–14
- Rajendran K, Harish CM (2000) Mechanism of triggered seismicity at Koyna: an evaluation based on relocated earthquakes. *Curr Sci* 79:358–363
- Rao NP, Shashidhar D (2016) Periodic variation of stress field in the Koyna-Warna reservoir triggered seismic zone inferred from focal mechanism studies. *Tectonophysics* 679:29–40
- Rao NP, Shashidhar D (2017) Earthquake focal mechanism studies in Koyna-Warna region in the last five decades. Current understanding on tectonics and seismogenesis. *J Geol Soc India* 90(6):684–691
- Rao VB, Murty BS, Murty AS (1969) Some geological and geophysical aspects of the Koyna (India) earthquake, December 1967. *Tectonophysics* 7(3):265–271
- Rao YB, Sreenivas B, Kumar TV, Khadke N, Krishna AK, Babu EVSSK (2017) Evidence for Neoproterozoic basement for the Deccan Volcanic flows around Koyna-Warna region, western India: Zircon U-Pb age and Hf-isotopic results. *J Geol Soc India* 90(6):752–760
- Rastogi BK, Talwani P (1980) Relocation of Koyna earthquakes. *Bull Seismol Soc Am* 70(5):1849–1868
- Rastogi BK, Chadha RK, Sarma CSP, Mandal P, Satyanarayana HVS, Raju IP, Kumar N, Satyamurthy C, Nageswara Rao A (1997) Seismicity at Warna reservoir (near Koyna) through 1995. *Bull Seismol Soc Am* 87(6):1484–1494
- Roeloffs E (1996) Poroelastic techniques in the study of earthquake-related hydrologic phenomena. In: *Advances in geophysics*, vol 37. Elsevier, Amsterdam, pp 135–195
- Roy S, Rao NP, Akkiraju VV, Goswami D, Sen M, Bansal BK, Nayak S (2013) Granitic basement below Deccan traps unearthed by drilling in the Koyna Seismic zone, western India. *J Geol Soc India* 81(2):289
- Sarkar D, Sain K (2017) Deep seismic sounding experiments in the Koyna RTS region: an overview of the results. *J Geol Soc India* 90(6):663–669
- Sarma SVS, Prasanta B, Patro K, Harinarayana T, Veeraswamy K, Sastry RS, Sarma MVC (2004) A magnetotelluric (MT) study across the Koyna seismic zone, western India: evidence for block structure. *Phys Earth Planet Inter* 142(1–2):23–36
- Sharma SR, Mall DM (1998) Geothermal and seismic evidence for the fluids in the crust beneath Koyna, India. *Curr Sci* 75:1070–1074
- Shashidhar D, Rao NP, Gupta H (2011) Waveform inversion of broadband data of local earthquakes in the Koyna-Warna region, western India. *Geophys J Int* 185(1):292–304
- Sibson RH (1985) A note on fault reactivation. *J Struct Geol* 7(6):751–754
- Singh C, Chadha RK (2010) Variations in the frequency–magnitude distribution of earthquakes with depth in the Koyna-Warna region, India. *J Asian Earth Sci* 39(4):331–334
- Singh DD, Rastogi BK, Gupta HK (1975) Surface-wave radiation pattern and source parameters of Koyna earthquake of December 10, 1967. *Bull Seismol Soc Am* 65(3):711–731
- Skempton AW (1961) Horizontal stresses in an overconsolidated Eocene clay. In: *Proceedings of the 5th international conference on soil mechanics*, vol 1, pp 351–357
- Srinagesh D, Rajagopala Sarma PVSS (2005) High precision earthquake locations in Koyna-Warna seismic zone reveal depth variation in brittle-ductile transition zone. *Geophys Res Lett* 32(8)
- Sykes LR (1970) Seismicity of the Indian Ocean and a possible nascent island arc between Ceylon and Australia. *J Geophys Res* 75(26):5041–5055
- Taghipour M, Ghafoori M, Lashkaripour GR, Moghaddas NH, Mola-ghab A (2019) Estimation of the current stress field and fault reactivation analysis in the Asmari reservoir, SW Iran. *Pet Sci* 16:513–526
- Talwani P (1995) Speculation on the causes of continuing seismicity near Koyna reservoir, India. In: *Induced seismicity*. Birkhäuser, Basel, pp 167–174
- Talwani P (1997a) Seismotectonics of the Koyna-Warna area, India. In: *Seismicity associated with mines, reservoirs and fluid injections*. Birkhäuser, Basel, pp 511–550
- Talwani P (1997b) On the nature of reservoir-induced seismicity. *Pure Appl Geophys* 150(3–4):473–492
- Talwani P (2000) Seismogenic properties of the crust inferred from recent studies of reservoir-induced seismicity-application to Koyna. *Curr Sci Bangalore* 79(9):1327–1333
- Talwani P, Kumarswamy SV, Sawalwede CB (1996) The re-evaluation of seismicity data in the Koyna-Warna area. Columbia University, South Carolina Publication, pp 1–109
- Talwani P, Cobb JS, Schaeffer MF (1999) In situ measurements of hydraulic properties of a shear zone in northwestern South Carolina. *J Geophys Res Solid Earth* 104(B7):14993–15003
- Tandon AN, Chaudhury HM (1968) Koyna earthquake of December 1967. Office of the Director General of Observatories
- Telesca L (2010) Analysis of the cross-correlation between seismicity and water level in the Koyna area of India. *Bull Seismol Soc Am* 100(5A):2317–2321

- Terzaghi K (1943) Theoretical soil mechanics. Wiley, New York, pp 11–15
- Tong H, Yin A (2011) Reactivation tendency analysis: a theory for predicting the temporal evolution of preexisting weakness under uniform stress state. *Tectonophysics* 503(3–4):195–200
- van Thienen-Visser K, Breunese JN (2015) Induced seismicity of the Groningen gas field: history and recent developments. *Lead Edge* 34(6):664–671
- Williams JN, Fagereng Å, Wedmore LN, Biggs J, Mphepo F, Dulanya Z, Mdala H, Blenkinsop T (2019) How do variably striking faults reactivate during rifting? Insights from southern Malawi. *Geochem Geophys Geosyst* 20:3588–3607
- Yadav A, Gahalaut K, Mallika K, Purnachandra Rao N (2015) Annual periodicity in the seismicity and water levels of the Koyna and Warna reservoirs, western India: a singular spectrum analysis. *Bull Seismol Soc Am* 105(1):464–472
- Yadav A, Bansal BK, Pandey AP (2016) Five decades of triggered earthquakes in Koyna-Warna Region, western India: a review. *Earth Sci Rev* 162:433–450
- Zoback MD, Zinke JC (2002) Production-induced normal faulting in the Valhall and Ekofisk oil fields. In: *The mechanism of induced seismicity*. Birkhäuser, Basel, pp 403–420



# Modelling reference evapotranspiration by combining neuro-fuzzy and evolutionary strategies

Meysam Alizamir<sup>1</sup> · Ozgur Kisi<sup>2</sup> · Rana Muhammad Adnan<sup>3</sup> · Alban Kuriqi<sup>4</sup>

Received: 5 January 2020 / Accepted: 15 May 2020 / Published online: 28 May 2020  
© Institute of Geophysics, Polish Academy of Sciences & Polish Academy of Sciences 2020

## Abstract

This study investigates the potential of two evolutionary neuro-fuzzy inference systems, adaptive neuro-fuzzy inference system (ANFIS) with particle swarm optimization (ANFIS–PSO) and genetic algorithm (ANFIS–GA), in modelling reference evapotranspiration ( $ET_0$ ). The hybrid models were tested using Nash–Sutcliffe efficiency, root mean square errors and determination coefficient ( $R^2$ ) statistics and compared with classical ANFIS, artificial neural networks (ANNs) and classification and regression tree (CART). Various combinations of monthly weather data of solar radiation, relative humidity, average air temperature and wind speed gotten from two stations, Antalya and Isparta, Turkey, were used as input parameters to the developed models to estimate  $ET_0$ . The recommended evolutionary neuro-fuzzy models produced better estimates compared to ANFIS, ANN and CART in modelling monthly  $ET_0$ . The ANFIS–PSO and/or ANFIS–GA improved the accuracy of ANFIS, ANN and CART by 40%, 32% and 66% for the Antalya and by 14%, 44% and 67% for the Isparta, respectively.

**Keywords** Reference evapotranspiration modelling · Evolutionary neuro-fuzzy inference systems · Particle swarm optimization · Genetic algorithm

## Introduction

Scarcity of water, increment in pumping costs, complications in water storage and delivery system are the main issues that emphasize on enhancement of the water application efficiency for the operation of large irrigation systems. Irrigation engineers and agricultural managers need to calculate crop water requirement accurately for utilizing the scarce water timely and efficiently. For the efficient water application, evapotranspiration (ET) has a crucial role due to help in the calculation of crop water requirements precisely. Therefore, an accurate estimation of ET is fundamental to

improve water application efficiency (Güven et al. 2008). The Food and Agriculture Organization (FAO) introduced the Penman–Monteith equation for modelling ET. This approach has become a commonly used method for calculating ET throughout the world (Allen et al. 2006). Several climatic inputs such as minimum, maximum and average temperature, wind speed, mean relative humidity and sunshine duration are required for ET estimation by the Penman–Monteith equation. These large numbers of climatic data are not always available or reliable. The influence of the mentioned climatic variables on ET makes it a complex nature (Hernandez et al. 2011), and therefore, forecasting ET is one of the most difficult tasks in water resource problems. In such a situation, soft computing (SC) methods that can accurately model complex behaviour between input and output emerge as a better alternative. In recent years, SC methods like ANNs, ANFIS and machine learning (ML) methods have applied for modelling different complex systems in the field of hydrology (Adnan et al. 2018, Adnan et al. 2019a, b; Nair et al. 2018; Muhammad Adnan et al. 2019; Majhi et al. 2019; Wu et al. 2020).

In the literature, ANNs and ANFIS models were applied successfully to predict evapotranspiration (Ladlani et al. 2012, 2014; Kisi et al. 2015; Wen et al. 2015; Luo et al.

✉ Meysam Alizamir  
meysamalizamir@gmail.com

<sup>1</sup> Department of Civil Engineering, Hamedan Branch, Islamic Azad University, Hamedan, Iran

<sup>2</sup> Faculty of Natural Sciences and Engineering, Ilia State University, Tbilisi, Georgia

<sup>3</sup> State Key Laboratory of Hydrology-Water Resources and Hydraulic Engineering, Hohai University, Nanjing 210098, China

<sup>4</sup> CERIS, Instituto Superior Técnico, Universidade de Lisboa, Lisbon, Portugal

2015; Keshtegar et al. 2018; Abrishami et al. 2019; Walls et al. 2020). Ladlani et al. (2012) compared two ANN models, namely generalized regression artificial neural network (GR-ANN) method and radial basis artificial neural network (RB-ANN), for modelling ET using climatic parameters from Dar El Beida, Algeria. As climatic data, the authors used the data of sunshine duration, average relative humidity, average wind speed, maximum, minimum and average air temperature. They found that the GR-ANN performed better than the RB-ANN in predicting ET. Ladlani et al. (2014) checked the potential of ANFIS and multiple linear regression (MLR) models for forecasting daily  $ET_0$  in the Mediterranean region of Algiers, Algeria. Results obtained from the investigation demonstrated that the ANFIS had better performance compared to the MLR models.

Kisi et al. (2015) compared four soft computing models: (1) MLP-ANN, (2) ANFIS-GP, (3) ANFIS-SC and (4) gene expression programming (GEP) models for predicting monthly ET using the data of 50 climatic stations in Iran. From the obtained results, the authors found the ANFIS-GP as an optimal model. Wen et al. (2015) investigated the prediction accuracy of ANN and empirical methods in comparison with a machine learning method, namely support vector machine (SVM). The selected models were used to predict ET of the arid region of Ejina basin, China, using the minimum temperature and maximum temperature as inputs. Luo et al. (2015) compared four ANN models: (i) multilayer perceptron artificial neural network (MLP-ANN), (2) generalized feed-forward artificial neural network (GFF-ANN), (3) probabilistic neural network (P-ANN) and (4) linear regression artificial neural network (LR-ANN) models for predicting evapotranspiration of Gaoyou climatic station of Jiangsu province in China. The results of this study proved that ANNs can be effectively employed as a reliable ET modelling tool. Keshtegar et al. (2018) applied the ANFIS (ANFIS-FCM) with ANN and M5 model tree models to predict the evapotranspiration of three stations of the Central Anatolian Region of Turkey. They divided data into different training–testing subsets to check ANFIS accuracy for each. They found that the ANFIS model with different subsets performed better than the M5 and ANN models. Abrishami et al. (2019) used the ANN models to predict the ET of Nissouri Creek in Oxford County, Canada. They used two types of activation functions including rectified linear unit (ReLU) and sigmoid. Results showed that ReLU performed better than sigmoid activation function. Walls et al. (2020) applied different ANN structures for modelling the ET of wheat and maize crops and found ANN models suitable for predicting ET of both crops. ANN, ANFIS and ML models have also been successfully used in modelling different hydrological time series due to their ability to capture nonlinear behaviour (Adnan et al. 2017; Kisi et al. 2018; Yuan et al. 2018).

In the recent past years, the literature study has exposed that the hybrid soft computing models provide better ET prediction accuracy in comparison with stand-alone soft computing methods. The primary consideration of the researchers is towards combining several novel heuristic search algorithms with soft computing methods for optimizing their control parameters and enhancement of their forecasting accuracy. Patil and Deka (2017) applied the hybrid of wavelet transform with ANN and ANFIS methods for the modelling of evapotranspiration in the arid regions of India. The results confirmed that the hybrid models had better performance than the stand-alone soft computing models in predicting ET. Araghi et al. (2018) also demonstrated the benefits of WT (wavelet transform) combined with the ANFIS (WT-ANFIS), ANN (WT-ANN) and MLR (WT-MLR) models for ET forecasting of three climatic stations chosen from three different climates of Iran. Using daily weather data of selected stations, the authors found that the WT-ANN outperformed the other wavelet-based hybrid models (i.e. WT-ANFIS and WT-MLR). Gocić et al. (2015) combined the firefly algorithm with SVM (SVM-FFA) for predicting ET in Serbia. The authors compared the proposed SVM-FFA model with WT-SVM, SVM and ANN. They found that the SVM-FFA and WT-SVM models provided better prediction results in comparison with stand-alone ANN and SVM computational methods. Shamshirband et al. (2016) applied a novel heuristic method called cuckoo search algorithm (CSA) for optimizing the ANN and ANFIS methods in estimation of ET at 12 climatic stations in Serbia. The prediction results of designed hybrid methods (ANN-CSA and ANFIS-CSA) are compared with stand-alone ANN and ANFIS models. Also, the authors compared the proposed methods with the Hargreaves and Priestley–Taylor empirical models.

Available literature indicates that hybrid heuristic soft computing methods generally provided better prediction accuracy compared to stand-alone soft computing models. The literature surveys point out that the application of new hybrid soft computing methods is vital to improve prediction accuracy and minimize the method's error. For this reason, evolutionary neuro-fuzzy systems are proposed in this research for an effective evapotranspiration modelling. Genetic algorithm (GA) and particle swarm optimization (PSO) heuristic algorithms are used to optimize the parameters of ANFIS models and to develop hybrid soft computing methods, ANFIS-PSO and ANFIS-GA. Also, ET modelling using classification and regression tree (CART) model is very scarce, and this study looks to be the first that compares the accuracy of CART with the ANFIS-PSO, ANFIS-GA, ANFIS and ANN models in ET prediction.

## Materials and methods

### Used data

The study uses monthly weather data, solar radiation, relative humidity, air temperature and wind speed, from two automated climatic stations, Antalya (long. of 30°44′00″E, lat. of 36°42′00″N and altitude of 64) and Isparta (long. of 30°34′00″E, lat. of 37°47′00″N and altitude of 997) operated by the TMO (Turkish Meteorological Organization). The study area and stations' location are illustrated in Fig. 1. The stations are situated in the Mediterranean region having a Mediterranean climate (dry summers and mellow to cold, wet winters). The temperature in winter has its highest value as 24 °C, and in summer season, it can increase to 40 °C.

In the study, data (25-year monthly values for the period of 1982–2006) were divided into two parts as training (80% of the aggregate data) and testing (20% remaining part). The brief statistical properties of the used data are summed up in Table 1. It is evident from the average statistics that the Antalya has a higher temperature, solar radiation, wind speed and reference evapotranspiration compared to Isparta.

### Used methods

#### Adaptive neuro-fuzzy inference system (ANFIS)

The ANFIS interface represents a multilayer model initially proposed by Jang (1993) that trains input and output variables and affords estimations agreement between input and output in the most efficient way. There are several fuzzy interfaces system (FIS) reported in the literature, which has different performance and as results in significant differences

in the results among them. The FIS is categorized into three main groups: Mamdani's interface system (Mamdani and Assilian 1975), which consists of a system that considered inputs and outputs as a fuzzy set. This system is the most often applied; Tsukamoto's system (Tsukamoto 1979), which is not very commonly used; finally, Sugeno's FIS, which considers the input data as a fuzzy set, while the outputs as a constant coefficient of a linear function (Takagi and Sugeno 1985). The fact of being compact and very efficient in terms of computational time makes the Sugeno's system very commonly used also (Nourani et al. 2014; Zhu et al. 2019; Adnan et al. 2019c; Alizamir et al. 2020a). ANFIS applied in this study consists of a network structure which uses Sugeno inference system (S-FIS) and supported from the artificial neural network (ANN) in the training phase of the input data (Fig. 2).

The ANFIS interface is composed of several nodes connected through directional links. Indeed, the combination of the fuzzy-based rules systems with the high performance regarding the learning capability of the ANN has made the ANFIS interface more robust and popular in modelling different problems (Tabari et al. 2012). ANFIS is more commonly used in solving complicated problems characterized by significantly high nonlinearity (Rezazakemi et al. 2017). Training of the data sets is done based on the fundamental learning rule backpropagation approach, which tends to minimize the error computation of the input data set (Cobaner 2011). In addition to the binary variables, a set of linguistic variables were used to design the fuzzy system. Afterwards, several IF/THEN rules were used to characterize the relationship between fuzzy variables (Nourani et al. 2014). In the case of the Sugeno's system, which is the system used in this study, the conditional rules IF/THEN can be expressed as follows (Sayed et al. 2003):

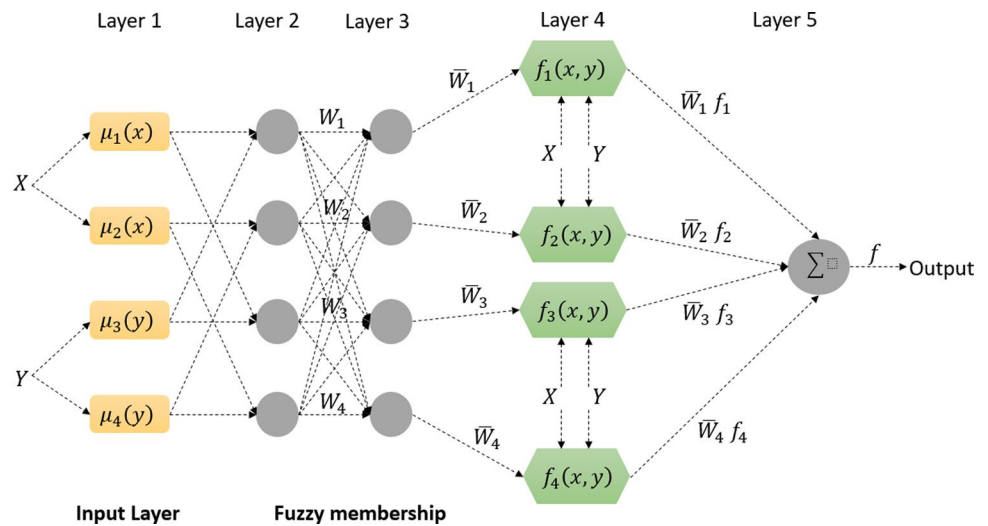
**Fig. 1** The study area and stations' location (adapted from d-maps.com)



**Table 1** Brief statistics for the climatic data of Antalya and Isparta stations

Station	Data set	Data set	Unit	Avr.	Min.	Max.	SD	Skewness
Antalya	Training data	$T$	°C	19.52	7.3	32.25	7.33	0.03
		SR	cal/cm <sup>2</sup>	412	120	679.2	154	-0.09
		RH	%	57.0	47.5	68.5	3.81	0.25
		WS	m/s	2.64	0.9	4.9	0.69	0.008
		ET <sub>0</sub>	mm day <sup>-1</sup>	5.64	1.16	10.4	2.1	0.16
	Testing data	$T$	°C	20.1	9.7	31.85	7.27	0.1
		SR	cal/cm <sup>2</sup>	361	1268	595.6	145	-0.08
		RH	%	52.9	45.5	67	4.16	0.93
Isparta	Training data	$T$	°C	12.3	-2.3	25	7.71	-0.12
		SR	cal/cm <sup>2</sup>	318	112	657	117	0.02
		RH	%	60.0	46	72	5.08	-0.31
		WS	m/s	1.84	0.6	3.6	0.5	0.42
		ET <sub>0</sub>	mm day <sup>-1</sup>	3.53	0.69	6.79	1.51	-0.02
	Testing data	$T$	°C	12.6	-0.9	25.2	8.04	-0.05
		SR	cal/cm <sup>2</sup>	355	148	638	141	0.21
		RH	%	63.4	52.5	72.5	2.57	0.006
		WS	m/s	1.43	0.8	2.5	0.42	0.68
		ET <sub>0</sub>	mm day <sup>-1</sup>	3.43	1.03	6.43	1.54	0.11

**Fig. 2** The fundamental structure of the ANFIS interface



Rule1 : if  $x$  is  $A_1$  and  $y$  is  $B_1$ , then  $z_1 = p_1x + q_1y + r_1$  (1)

Rule2: if  $x$  is  $A_2$  and  $y$  is  $B_2$ , then  $z_2 = p_2x + q_2y + r_2$  (2)

where  $A_1$  and  $B_1$  represent the fuzzy sets in the originator, and  $p_i$ ,  $q_i$  and  $r_i$  are the design parameters defined during the training process of the data set. As shown in Fig. 2, the architecture of the ANFIS was designed considering five layers, and detailed explanation for each layer and the equations used can be found in the literature (Tabari et al. 2012). The hybrid learning algorithm used in the ANFIS architecture

applies a combination of gradient descent, in order to identify the proposition parameters, while the least-squares method is applied to allocate the linear consequent parameters. The training algorithm makes the ANFIS outputs with the lowest error (Jang 1993; Nourani et al. 2014).

**ANFIS-PSO** In this ANFIS model, a particle swarm optimization (PSO) was used. This optimization algorithm is very efficient in case of discrete data type (Nourani et al. 2014). This combination may be considered as a surrogate approach. So, after determining the design variables, the objective function and constraints, ANFIS was mainly used

to search the space, while PSO approach as an optimization algorithm can be employed to establish the efficient way to find the best salutation (Kennedy and Eberhart 1995). PSO is a stochastic optimization method which consists of selecting a specific population or particles in given space completely in a random way while subsequently looking for the optimal solution (Rezakazemi et al. 2017). There are several applications of ANFIS–PSO. Bassar et al. (2015) used ANFIS–PSO to predict the optimal parameters to mitigate scouring depth in existing spur dykes, while Djavahreshkian and Esmaeili (2014) applied ANFIS–PSO to optimize the operation of the submerged hydrofoil. ANFIS–PSO interface is also used to solve nonlinear problems related to the nanomaterial’s components (Rezakazemi et al. 2017).

**ANFIS–GA** In ANFIS–GA model, genetic optimization algorithm (GA) is incorporated. The GA consists of the inset of chromosome combinations, which evaluates the results obtained in each computational step by seeking the optimal solution possible (Termeh et al. 2018). Differently, from PSO, GA can provide relatively large solution spaces since it utilizes a probabilistic transition and not deterministic rules (Rezakazemi et al. 2017). The GA interface uses variables that represent real values or binary coding. The GA optimization procedure is associated with several processes as follows: population initialization, selection, crossover and mutation (Rezakazemi et al. 2017). GA has become a prevalent optimization method in different areas. Termeh et al. (2018) applied GA in flood susceptibility mapping; they found that this algorithm among the other advantages reveals high accuracy. Rezakazemi et al. (2017) used GA to assess the hydrogen mixed matrix membrane considering several operating conditions. While Khosravi et al. (2018)

applied GA to predict potential solar radiation to support solar-based energy systems, all the studies above found that the GA interface poses the ability to provide efficient computation time and high accuracy.

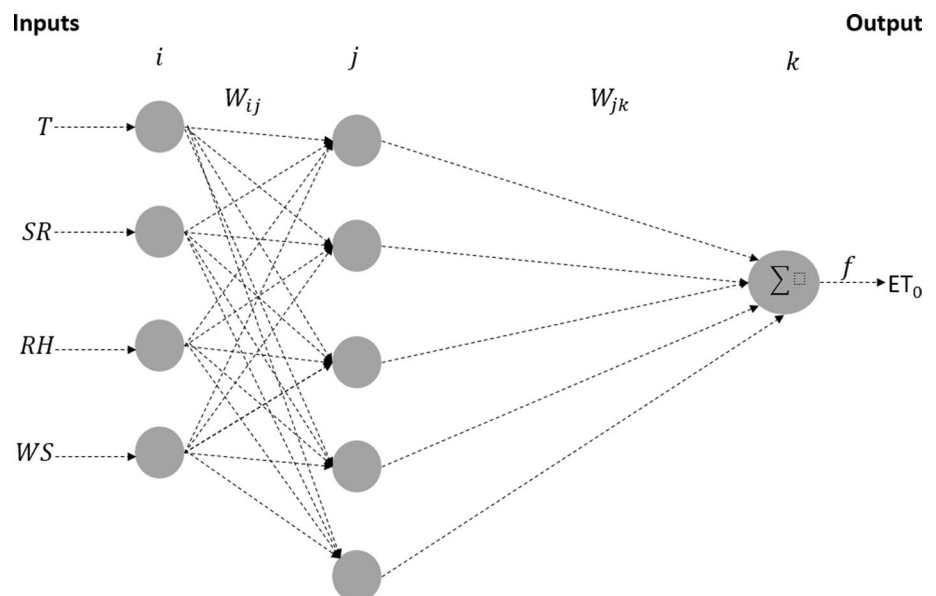
### Artificial neural network (ANN)

Artificial neural network (ANN) consists of imitating the biological nervous system, although much of the biological details are neglected. ANNs are composed of several massively processing elements organized in parallel systems connected by using variable weights. Each layer is connected to the other layers through interconnection weights,  $W$ . The methodology applied for tuning the weights based on backpropagation process (Rumelhart et al. 1986). The backpropagation network is by far the most commonly used paradigms in ANNs (Nourani et al. 2014; Kisi et al. 2017; Alizamir et al. 2018; Kisi and Alizamir 2018). The processing elements that composed the ANNs are called neurons. The basic structure of the ANN interface is shown in Fig. 3. The neural network layers  $i$ ,  $j$  and  $k$  are interconnected with weights  $W_{ij}$  and  $W_{jk}$  between layers of neurons. Further details and explanation about the training process of the input data may be found at Kisi and Öztürk (2007).

### Classification and regression tree (CART)

Classification and regression tree (CART) is based on a set of decision trees on the predictor variables which grew by repeatedly stratifying the data set into consecutively smaller subgroups (Breiman 1984). CART is a predictive tree model based on the recursive approach in data mining models that constructs the structure of the given data set which generates

**Fig. 3** The ANN interface used for the  $ET_0$  estimation

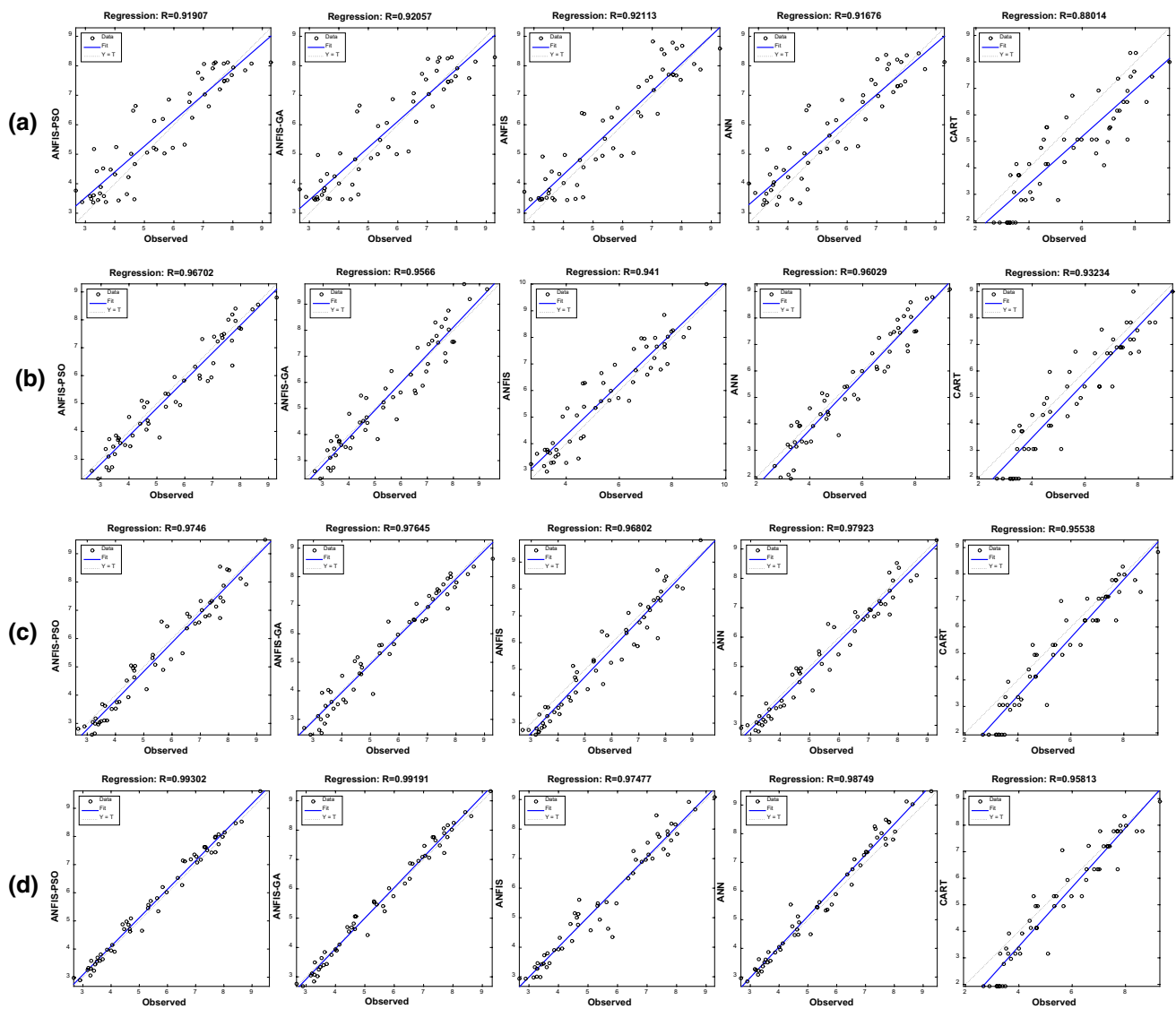




**Table 2** The test statistics of the ANFIS-PSO, ANFIS-GA, ANFIS, ANN and CART models in estimating ET<sub>0</sub> of Antalya Station

Input combination	RMSE					NSE					R <sup>2</sup>				
	ANFIS-PSO	ANFIS-GA	ANFIS	ANN	CART	ANFIS-PSO	ANFIS-GA	ANFIS	ANN	CART	ANFIS-PSO	ANFIS-GA	ANFIS	ANN	CART
(i)															
T	<b>0.742</b>	<b>0.730</b>	<b>0.762</b>	<b>0.758</b>	1.230	<b>0.831</b>	<b>0.837</b>	<b>0.822</b>	<b>0.824</b>	0.537	<b>0.844</b>	<b>0.847</b>	<b>0.848</b>	<b>0.84</b>	0.618
SR	1.023	0.952	1.033	1.474	<b>1.179</b>	0.679	0.722	0.673	0.335	<b>0.574</b>	0.819	0.835	0.819	0.793	<b>0.774</b>
RH	2.058	2.090	2.129	1.785	2.170	-0.295	-0.335	-0.386	0.025	-0.439	0.001	0.003	0.023	0.151	0.024
WS	1.829	1.837	1.848	1.835	1.858	-0.023	-0.031	-0.043	-0.029	-0.056	0.047	0.046	0.036	0.041	0.071
(ii)															
T, SR	0.780	0.718	0.871	1.717	0.904	0.814	0.842	0.768	0.098	0.750	0.889	0.887	0.853	0.149	0.865
T, RH	0.831	0.833	0.963	0.763	1.807	0.788	0.787	0.716	0.822	0.001	0.863	0.813	0.734	0.858	0.223
T, WS	0.639	0.617	<b>0.681</b>	0.652	0.887	0.875	0.883	<b>0.858</b>	0.869	0.759	0.905	0.898	<b>0.885</b>	0.901	0.792
SR, RH	<b>0.502</b>	<b>0.593</b>	0.779	<b>0.58</b>	<b>0.861</b>	<b>0.922</b>	<b>0.892</b>	0.814	<b>0.896</b>	<b>0.773</b>	<b>0.935</b>	<b>0.915</b>	0.906	<b>0.922</b>	<b>0.869</b>
SR, WS	0.898	0.880	0.776	1.424	1.109	0.753	0.763	0.815	0.379	0.623	0.833	0.856	0.872	0.802	0.808
(iii)															
T, SR, RH	0.588	<b>0.402</b>	0.908	0.437	0.868	0.894	<b>0.950</b>	0.747	0.941	0.769	0.910	<b>0.953</b>	0.857	0.954	0.875
T, SR, WS	<b>0.453</b>	0.441	<b>0.555</b>	<b>0.399</b>	<b>0.774</b>	<b>0.937</b>	0.940	<b>0.905</b>	<b>0.951</b>	<b>0.816</b>	<b>0.949</b>	0.957	<b>0.937</b>	<b>0.958</b>	<b>0.912</b>
(iv)															
T, SR, RH, WS	<b>0.248</b>	<b>0.253</b>	<b>0.424</b>	<b>0.371</b>	<b>0.746</b>	<b>0.981</b>	<b>0.980</b>	<b>0.945</b>	<b>0.957</b>	<b>0.829</b>	<b>0.986</b>	<b>0.983</b>	<b>0.950</b>	<b>0.975</b>	<b>0.918</b>

The bold numbers show the best model in each input combination



**Fig. 4** The observed and estimated  $ET_0$  values by the best models in the test period of Antalya Station: **a** input combination (i), **b** input combination (ii), **c** input combination (iii) and **d** input combination (iv)

decision rules for predicting a categorical variable (Choubin et al. 2018; Kisi et al. 2020; Alizamir et al. 2020a, b). Considering the principle of homogenization or less variability among the nodes, the splitting procedure of the variables is made until the best split is reached (Breiman 1984).

CART algorithm has also become commonly applied in different fields. Choubin et al. (2018) have applied CART to predict sediment transport in alpine rivers; they found that CART has relatively high accuracy. Ebrahimi and Azadbakht (2019) have applied CART to predict land surface temperature over several different areas. Also, Juntakut et al. (2019) have used CART to predict the long-term contamination of the groundwater in Nebraska State. They concluded that CART was capable of differentiating the weight of several physical factors in the water contamination.

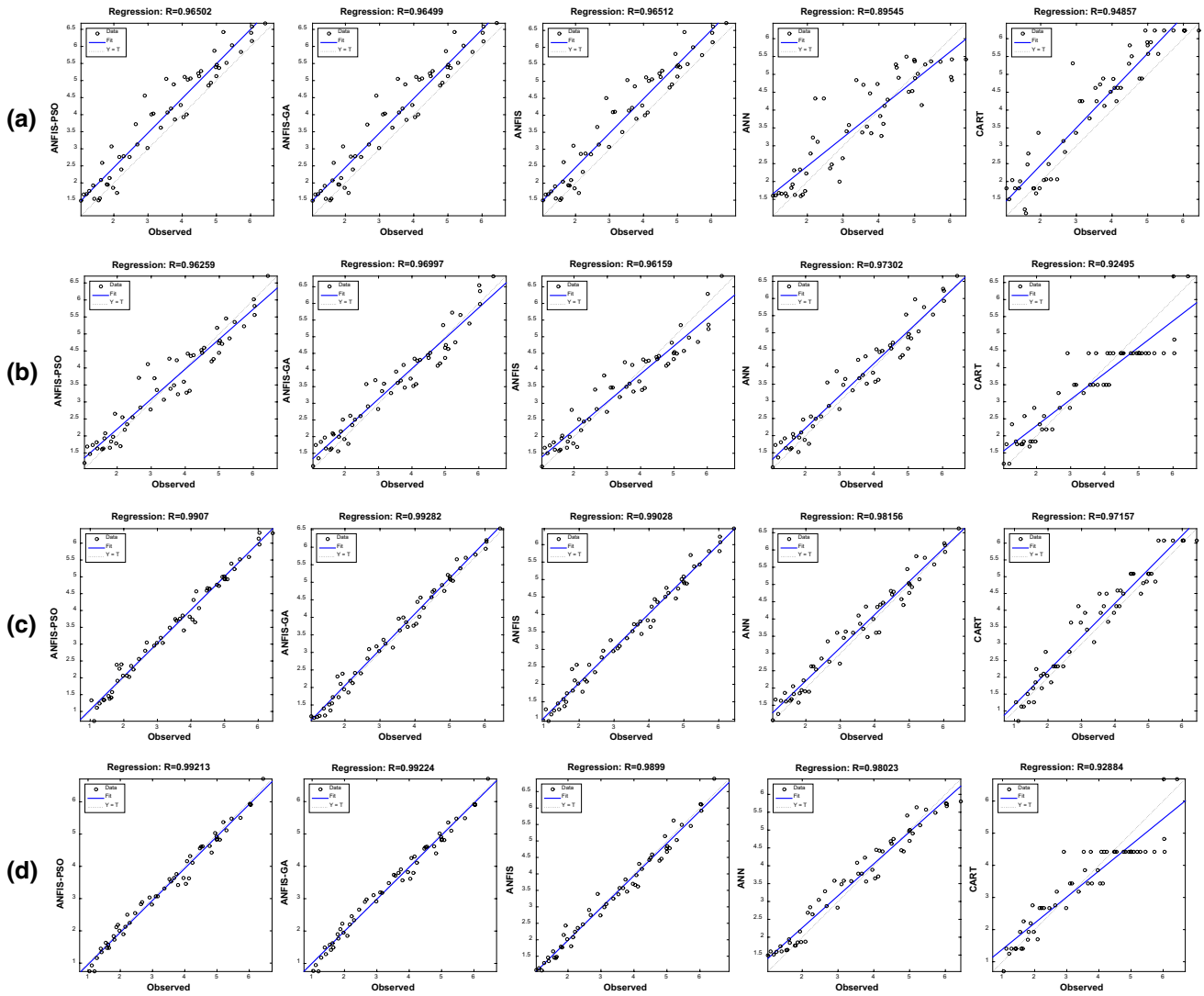
## Application and results

The ability of two evolutionary neuro-fuzzy systems: ANFIS-PSO and ANFIS-GA, are investigated in modeling reference evapotranspiration ( $ET_0$ ) using various input combinations of climatic data and compared with the classic ANFIS, ANN and CART methods. For the control parameters, different values were tried for each method. For the ANFIS-PSO, 500 iterations were used, and population, inertia weight, personal learning coefficient and global learning coefficient were set to 45, 1, 1 and 2, respectively. For the ANFIS-GA, number of iterations, population, mutation and crossover percentages were set to 400, 55, 0.7 and 0.4, respectively. For the ANFIS, subtractive clustering with 150 iterations and 0.35 radii was used. For the ANN, Bayesian

**Table 3** The test statistics of the ANFIS-PSO, ANFIS-GA, ANFIS, ANN and CART models in estimating ET<sub>0</sub> of Isparta Station

Input combination	RMSE					NSE					R <sup>2</sup>				
	ANFIS-PSO	ANFIS-GA	ANFIS	ANN	CART	ANFIS-PSO	ANFIS-GA	ANFIS	ANN	CART	ANFIS-PSO	ANFIS-GA	ANFIS	ANN	CART
(i)															
T	0.683	0.686	0.694	<b>0.697</b>	0.829	0.800	0.799	0.794	<b>0.792</b>	0.706	0.808	0.805	0.799	<b>0.801</b>	0.711
SR	<b>0.632</b>	<b>0.631</b>	<b>0.638</b>	1.306	<b>0.748</b>	<b>0.829</b>	<b>0.829</b>	<b>0.825</b>	0.271	<b>0.761</b>	<b>0.931</b>	<b>0.931</b>	<b>0.931</b>	0.850	<b>0.899</b>
RH	1.494	1.480	1.502	1.514	1.456	0.047	0.065	0.036	0.021	0.095	0.271	0.273	0.268	0.270	0.262
WS	1.542	1.524	1.546	1.557	1.608	-0.015	0.007	-0.02	-0.034	-0.104	0.001	0.014	0.001	0.006	0.002
(ii)															
T, SR	0.492	0.515	0.492	1.317	0.627	0.896	0.886	0.896	0.259	0.832	0.961	0.960	0.964	0.724	0.906
T, RH	0.750	0.706	0.790	0.702	0.777	0.760	0.786	0.733	0.789	0.742	0.800	0.811	0.791	0.816	0.783
T, WS	0.678	0.669	0.774	0.675	0.804	0.803	0.809	0.743	0.805	0.723	0.852	0.844	0.786	0.848	0.745
SR, RH	<b>0.422</b>	<b>0.386</b>	<b>0.441</b>	<b>0.380</b>	<b>0.600</b>	<b>0.923</b>	<b>0.936</b>	<b>0.916</b>	<b>0.938</b>	<b>0.846</b>	<b>0.926</b>	<b>0.940</b>	<b>0.924</b>	<b>0.946</b>	<b>0.855</b>
SR, WS	0.525	0.511	0.545	1.472	0.651	0.882	0.888	0.873	0.074	0.818	0.912	0.911	0.899	0.811	0.844
(iii)															
T, SR, RH	0.333	0.338	0.319	<b>0.329</b>	0.599	0.952	0.951	0.956	<b>0.953</b>	0.846	0.953	0.956	0.959	<b>0.963</b>	0.859
T, SR, WS	<b>0.210</b>	<b>0.212</b>	<b>0.215</b>	1.302	<b>0.430</b>	<b>0.981</b>	<b>0.980</b>	<b>0.980</b>	0.270	<b>0.921</b>	<b>0.981</b>	<b>0.985</b>	<b>0.980</b>	0.863	<b>0.944</b>
(iv)															
T, SR, RH, WS	<b>0.201</b>	<b>0.191</b>	<b>0.221</b>	<b>0.339</b>	<b>0.577</b>	<b>0.982</b>	<b>0.984</b>	<b>0.979</b>	<b>0.950</b>	<b>0.857</b>	<b>0.984</b>	<b>0.984</b>	<b>0.979</b>	<b>0.960</b>	<b>0.862</b>

The bold numbers show the best model in each input combination



**Fig. 5** The observed and estimated  $ET_0$  values by the best models in the test period of Isparta Station: **a** input combination (i), **b** input combination (ii), **c** input combination (iii) and **d** input combination (iv)

regulation was used, and the optimal number of neurons in the hidden layer (HL) was found to be 10. The following evaluation metrics are used to select the best models:

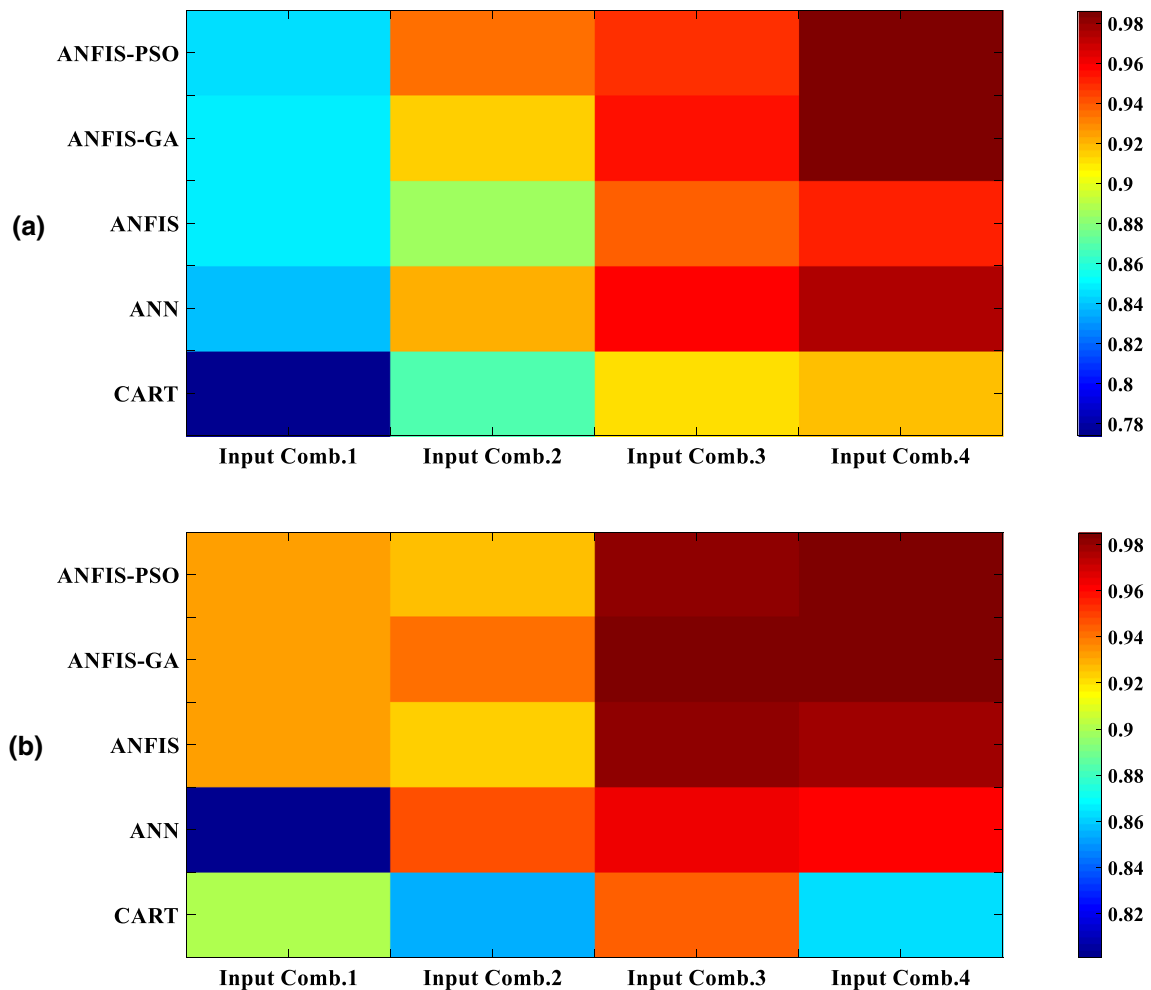
$$\text{Root mean square error (RMSE)} = \sqrt{\frac{\sum_{i=1}^N (ET_{im} - ET_{ip})^2}{N}} \tag{3}$$

$$\text{Nash – Sutcliffe efficiency (NSE)} = 1 - \frac{\sum_{i=1}^N (ET_{im} - ET_{ip})^2}{\sum_{i=1}^N (ET_{im} - \overline{ET}_m)^2} \tag{4}$$

where  $N$  = number of data,  $\overline{ET}_m$  = mean FAO 56 PM  $ET_0$ ,  $ET_{ip}$  = predicted  $ET_0$ , and  $ET_{im}$  = FAO 56 PM  $ET_0$ .

Table 2 compares the test statistics of the ANFIS–PSO, ANFIS–GA, ANFIS, ANN and CART models for different

input combinations of Antalya Station. Among the one input combinations,  $T$  variable provided the best statistics for all methods. Out of two-input models, the model with SR and RH inputs had the lowest RMSE and the highest NSE and  $R^2$  for the ANFIS–PSO, ANFIS–GA, ANN and CART methods. Three-input ANFIS–PSO, ANFIS, ANN and CART models with  $T$ , SR and WS variables performed better than the corresponding models with  $T$ , SR and RH variables. It is apparent from Table 2 that the models with whole input variables ( $T$ , SR, RH and WS) outperformed the other models for all methods. The ANFIS–PSO and ANFIS–GA with full climatic inputs have almost the same accuracy, and they have better statistics than the other models. The relative RMSE differences between the ANFIS–PSO and ANFIS, ANN, CART are 40%, 32% and 66%, respectively.



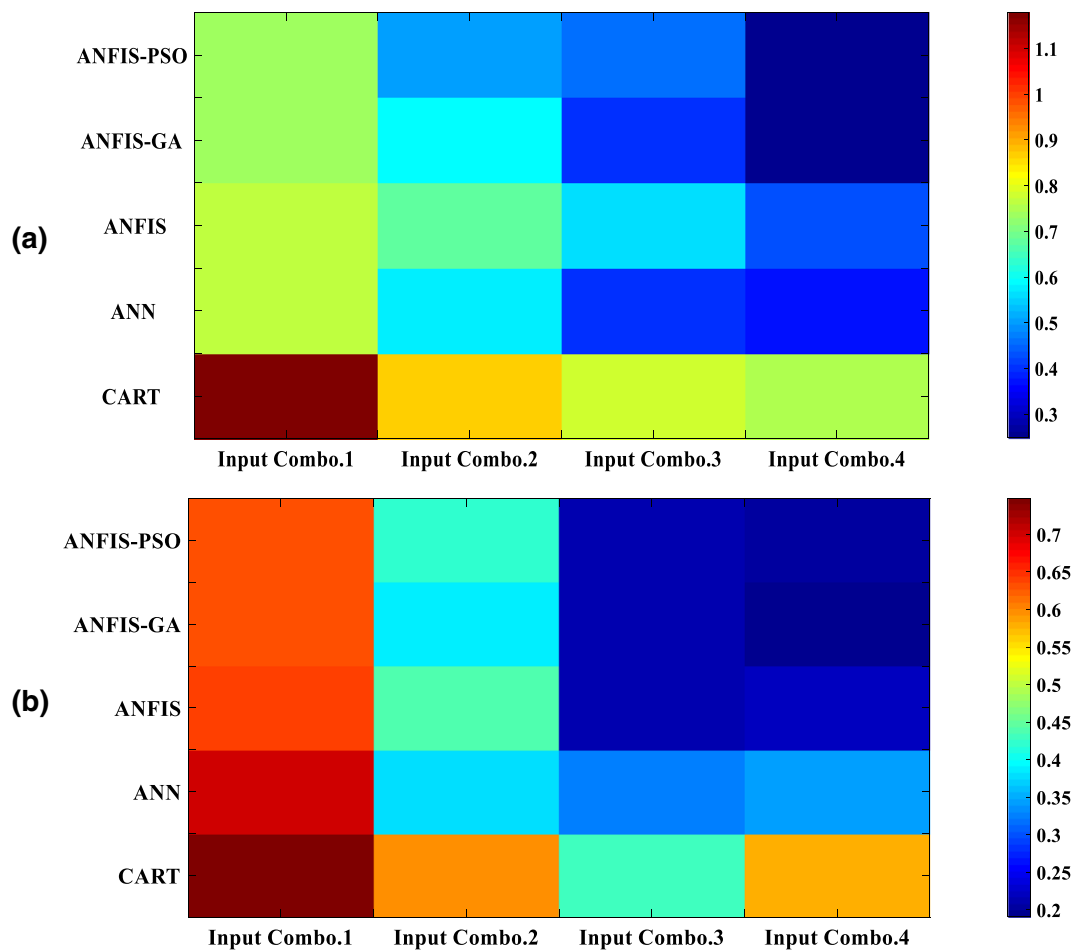
**Fig. 6**  $R^2$  values of the best models for different input combinations: **a** Antalya Station, **b** Isparta Station

Figure 4 illustrates the observed and estimated  $ET_0$  values by the best models in the test period for the Antalya Station. It is clearly observed that the ANFIS–PSO and/or ANFIS–GA models generally have less scattered estimates compared to other models. It is also apparent from the scatter graphs that all the methods produce less scattered estimates by increasing the number of input variables.

Test results of the employed methods are summed up in Table 3 for estimating  $ET_0$  of Isparta Station utilizing various climatic input variables. In this station, the models with SR variable have the best statistics among the one input combinations. Similar to the Antalya Station, here also the SR, RH and T, SR, WS combinations generally provided the most accurate estimates for two- and three-input models, respectively. Among all input combinations, the models with full climatic input variables performed the best. The best ANFIS–GA model outperformed the ANFIS–PSO, ANFIS, ANN and CART with respect to RMSE, NSE and  $R^2$ . The

relative RMSE differences between the ANFIS–GA and ANFIS–PSO, ANFIS, ANN, CART are 5%, 14%, 44% and 67%, respectively. It is clear from Tables 2 and 3 that the evolutionary algorithms, PSO and GA, improve the classical ANFIS model in both stations, improvement in RMSE by about 40% and 14% for the Antalya and Isparta stations. In Isparta Station, SR seems to be more effective on  $ET_0$  compared to Antalya. The RH and WS variables produce worse results compared to Isparta. One reason for this may be the fact that these variables have higher skewed distribution in Antalya than the Isparta (see skewness values of the SR and RH data in Table 1).

The test results of the employed models are graphically compared in Fig. 5 for the Isparta Station. Here also, the better estimates are obtained by increasing input numbers, and the models with full inputs (T, SR, RH and WS) have the best estimates among the input combinations tried. Both ANFIS–PSO and ANFIS–GA have less scattered



**Fig. 7** RMSE values of the best models for different input combinations: **a** Antalya Station, **b** Isparta Station

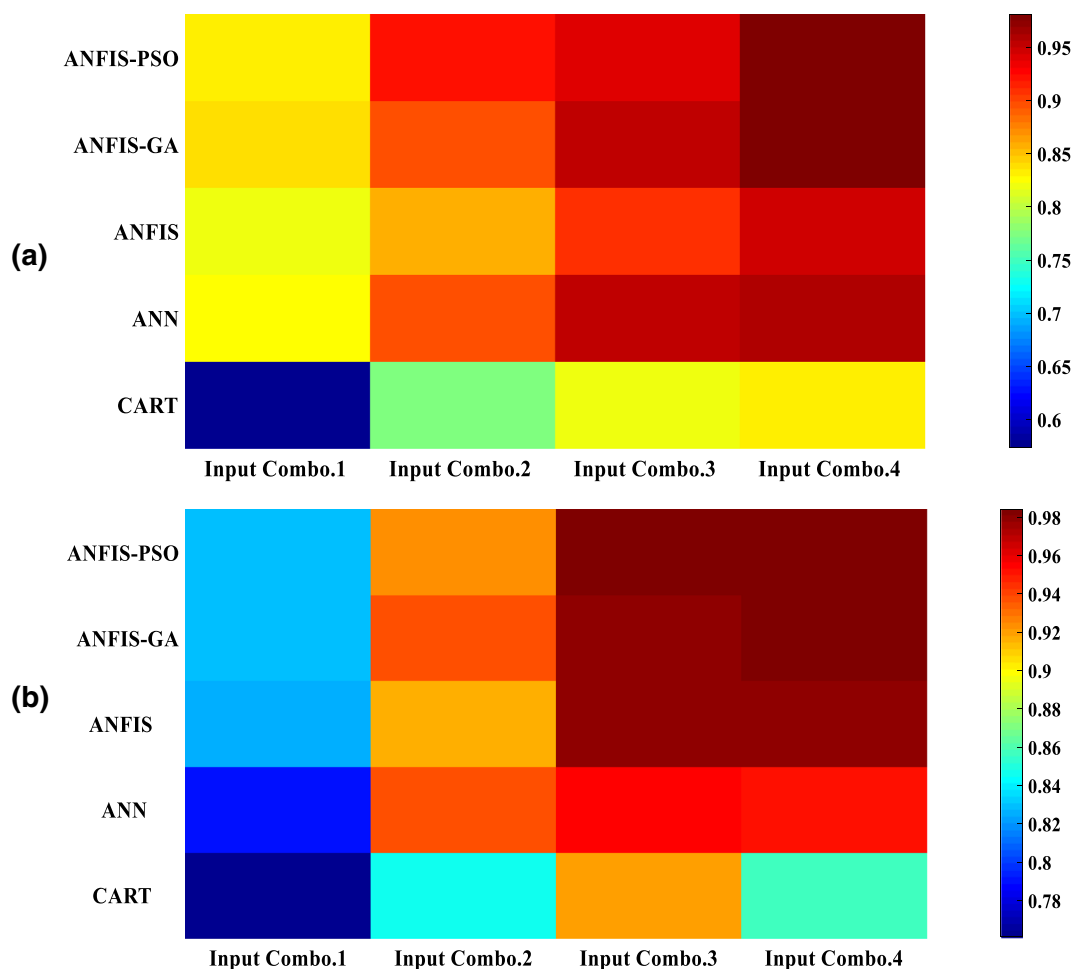
estimates than the ANFIS, ANN and CART models. CART model has the worst estimates among the models applied. Figures 6, 7 and 8 compare the  $R^2$ , RMSE and NSE values of the five optimal models with different input combinations for the Antalya and Isparta stations, respectively. It is clearly observed that the ANFIS–PSO and/or ANFIS–GA generally have the highest  $R^2$  and NSE and the lowest RMSE compared to other three methods.

Overall, the ANFIS–PSO and ANFIS–GA models perform superior to the other models in estimating monthly  $ET_0$ . PSO and GA are heuristic methods and have some advantages compared to classical training algorithms such as gradient descent and least square. These belong to a class of search methods so that they have a notable balance between exploitation of the optimal solutions and reconnaissance of the search space. Stochastic search and directed search are combined in such methods. Therefore, they are more robust compared to directed search techniques and capable of finding global optimum without local optima problem (Mantoglou et al. 2004; Karterakis et al. 2007).

## Conclusion

The accuracy of two evolutionary neuro-fuzzy methods was investigated in the presented study in modelling reference evapotranspiration. Their results were compared with the classic ANFIS, ANN and CART models. Various input combinations of climatic data obtained from two stations; Turkey were utilized for the employed models. Evolutionary ANFIS–PSO and/or ANFIS–GA produced better  $ET_0$  estimates than the ANFIS, ANN and CART models with the relative RMSE differences of 40%, 32% and 66% for one station (Antalya) and 14%, 44% and 67% for the other station (Isparta), respectively.

Comparison of various climatic inputs revealed that the estimation accuracy of the applied models increases by including more input variables and four inputs (average temperature, solar radiation, relative humidity and wind speed) produced the best estimates for each method. The comparison also indicated that solar radiation has more influence on  $ET_0$  in Isparta, while including relative humidity and wind speed in inputs makes models less accurate in Antalya.



**Fig. 8** NSE values of the best models for different input combinations: **a** Antalya Station, **b** Isparta Station

**Acknowledgement** The authors would like to thank the anonymous reviewers for their valuable comments and suggestions to improve this manuscript further. Alban Kuriqi was supported by a Ph.D. scholarship granted by Fundação para a Ciência e a Tecnologia, I.P. (FCT), Portugal, under the Ph.D. Programme FLUVIO–River Restoration and Management, grant number: PD/BD/114558/2016.

### Compliance with ethical standards

**Conflict of interest** There is not any conflict of interest in this study.

### References

- Abrishami N, Sepaskhah AR, Shahrokhnia MH (2019) Estimating wheat and maize daily evapotranspiration using artificial neural network. *Theoret Appl Climatol* 135(3–4):945–958
- Adnan RM, Yuan X, Kisi O, Anan R (2017) Improving accuracy of river flow forecasting using LSSVR with gravitational search algorithm. *Adv Meteorol* 2017:23
- Adnan RM, Yuan X, Kisi O, Adnan M, Mehmood A (2018) Stream flow forecasting of poorly gauged mountainous watershed by least square support vector machine, fuzzy genetic algorithm and M5 model tree using climatic data from nearby station. *Water Resour Manag* 32(14):4469–4486
- Adnan RM, Liang Z, Yuan X, Kisi O, Akhlaq M, Li B (2019a) Comparison of LSSVR, M5RT, NF-GP, and NF-SC models for predictions of hourly wind speed and wind power based on cross-validation. *Energies* 12(2):329
- Adnan RM, Liang Z, Trajkovic S, Zounemat-Kermani M, Li B, Kisi O (2019b) Daily streamflow prediction using optimally pruned extreme learning machine. *J Hydrol* 577:123981
- Adnan RM, Malik A, Kumar A, Parmar KS, Kisi O (2019c) Pan evaporation modeling by three different neuro-fuzzy intelligent systems using climatic inputs. *Arab J Geosci* 12(20):606
- Alizamir M, Kisi O, Zounemat-Kermani M (2018) Modelling long-term groundwater fluctuations by extreme learning machine using hydro-climatic data. *Hydrol Sci J* 63(1):63–73
- Alizamir M, Kim S, Kisi O, Zounemat-Kermani M (2020a) A comparative study of several machine learning based non-linear regression methods in estimating solar radiation: case studies of the USA and Turkey regions. *Energy* 197:117239
- Alizamir M, Kim S, Kisi O, Zounemat-Kermani M (2020b) Deep echo state network: a novel machine learning approach to model dew point temperature using meteorological variables. *Hydrol Sci J* 65(7):1173–1190

- Allen RG, Pruitt WO, Wright JL, Howell TA, Ventura F, Snyder R, Itenfisu D, Stedudo P, Berengena J, Yrisarry JB, Smith M, Raes D, Perrier A, Alves I, Walter I, Elliot R (2006) A recommendation on standardized surface resistance for hourly calculation of reference ETo by the FAO 56 Penman–Monteith method. *Agric Water Manag* 81:1–22
- Araghi A, Adamowski J, Martinez CJ (2018) Comparison of wavelet-based hybrid models for the estimation of daily reference evapotranspiration in different climates. *J Water Clim Change*. <https://doi.org/10.2166/wcc.2018.113>
- Basser H, Karami H, Shamshirband S, Akib S, Amirmojahedi M, Ahmad R, Javidnia H (2015) Hybrid ANFIS–PSO approach for predicting optimum parameters of a protective spur dike. *Appl Soft Comput* 30:642–649
- Breiman L (1984) *Classification and regression trees*. Chapman & Hall, London
- Choubin B, Darabi H, Rahmati O, Sajedi-Hosseini F, Kløve B (2018) River suspended sediment modelling using the CART model: a comparative study of machine learning techniques. *Sci Total Environ* 615:272–281
- Cobaner M (2011) Evapotranspiration estimation by two different neuro-fuzzy inference systems. *J Hydrol* 398(3–4):292–302
- Djavareshkian MH, Esmaili A (2014) Heuristic optimization of submerged hydrofoil using ANFIS–PSO. *Ocean Eng* 92:55–63
- Ebrahimi H, Azadbakht M (2019) Downscaling MODIS land surface temperature over a heterogeneous area: an investigation of machine learning techniques, feature selection, and impacts of mixed pixels. *Comput Geosci* 124:93–102
- Gocić M, Motamedi S, Shamshirband S, Petković D, Ch S, Hashim R, Arif M (2015) Soft computing approaches for forecasting reference evapotranspiration. *Comput Electron Agric* 113:164–173
- Guvan A, Aytek A, Yuce MI, Aksoy H (2008) Genetic programming-based empirical model for daily reference evapotranspiration estimation. *CLEAN Soil Air Water* 36(10–11):905–912
- Hernandez S, Morales L, Sallis P (2011) Estimation of reference evapotranspiration using limited climatic data and Bayesian model averaging. In: *Ems, 2011 UK Sim 5th European symposium on computer modeling and simulation*. pp 59–63
- Jang JS (1993) ANFIS: adaptive-network-based fuzzy inference system. *IEEE Trans Syst Man Cybern* 23(3):665–685
- Juntakut P, Snow DD, Haacker EM, Ray C (2019) The long term effect of agricultural, vadose zone and climatic factors on nitrate contamination in the Nebraska’s groundwater system. *J Contam Hydrol* 220:33–48
- Karterakis SM, Karatzas GP, Nikolos IK, Papadopoulou MP (2007) Application of linear programming and differential evolutionary optimization methodologies for the solution of coastal subsurface water management problems subject to environmental criteria. *J Hydrol* 342(3–4):270–282
- Kennedy J, Eberhart R (1995) Particle swarm optimization. *IEEE Int Conf 4:1942–1948*
- Keshtegar B, Kisi O, Arab HG, Zounemat-Kermani M (2018) Subset modeling basis ANFIS for prediction of the reference evapotranspiration. *Water Resour Manag* 32(3):1101–1116
- Khosravi A, Nunes RO, Assad MEH, Machado L (2018) Comparison of artificial intelligence methods in estimation of daily global solar radiation. *J Clean Prod* 194:342–358
- Kisi O, Alizamir M (2018) Modelling reference evapotranspiration using a new wavelet conjunction heuristic method: wavelet extreme learning machine vs wavelet neural networks. *Agric For Meteorol* 263:41–48
- Kişİ Ö, Öztürk Ö (2007) Adaptive neuro fuzzy computing technique for evapotranspiration estimation. *J Irrig Drain Eng* 133(4):368–379
- Kisi O, Sanikhani H, Zounemat-Kermani M, Niazi F (2015) Long-term monthly evapotranspiration modeling by several data-driven methods without climatic data. *Comput Electron Agric* 115:66–77
- Kisi O, Alizamir M et al (2017) Modeling groundwater fluctuations by three different evolutionary neural network techniques using hydroclimatic data. *Nat Hazards* 87(1):367–381
- Kisi O, Shiri J, Karimi S, Adnan RM (2018) Three different adaptive neuro fuzzy computing techniques for forecasting long-period daily streamflows. In: Roy SS et al (eds) *Big data in engineering applications*. Springer, Singapore, pp 303–321
- Kisi O, Alizamir M, Docheshmeh Gorgij A (2020) Dissolved oxygen prediction using a new ensemble method. *Environ Sci Pollut Res* 27:9589–9603
- Ladlani I, Houichi L, Djemili L, Heddam S, Belouz K (2012) Modeling daily reference evapotranspiration (ET<sub>0</sub>) in the north of Algeria using generalized regression neural networks (GRNN) and radial basis function neural networks (RBFNN): a comparative study. *Meteorol Atmos Phys* 118(3–4):163–178
- Ladlani I, Houichi L, Djemili L, Heddam S, Belouz K (2014) Estimation of daily reference evapotranspiration (ET<sub>0</sub>) in the north of Algeria using adaptive neuro-fuzzy inference system (ANFIS) and multiple linear regression (MLR) models: a comparative study. *Arab J Sci Eng* 39(8):5959–5969
- Luo Y, Traore S, Lyu X, Wang W, Wang Y, Xie Y, Fipps G (2015) Medium range daily reference evapotranspiration forecasting by using ANN and public weather forecasts. *Water Resour Manag* 29(10):3863–3876
- Majhi B, Naidu D, Mishra AP, Satapathy SC (2019) Improved prediction of daily pan evaporation using deep-LSTM model. *Neural Comput Appl*. <https://doi.org/10.1007/s00521-019-04127-7>
- Mamdani EH, Assilian S (1975) An experiment in linguistic synthesis with a fuzzy logic controller. *Int J Man Mach Stud* 7(1):1–13
- Mantoglou A, Papantoniou M, Giannouloupoulos P (2004) Management of coastal aquifers based on nonlinear optimization and evolutionary algorithms. *J Hydrol* 297(1–4):209–228
- Muhammad Adnan R, Yuan X, Kisi O, Yuan Y, Tayyab M, Lei X (2019) Application of soft computing models in streamflow forecasting. In: *Proceedings of the institution of civil engineers-water management*. Thomas Telford Ltd., vol 172(3), pp 123–134
- Nair A, Singh G, Mohanty UC (2018) Prediction of monthly summer monsoon rainfall using global climate models through artificial neural network technique. *Pure appl Geophys* 175(1):403–419
- Nourani V, Baghanam AH, Adamowski J, Kisi O (2014) Applications of hybrid wavelet–artificial intelligence models in hydrology: a review. *J Hydrol* 514:358–377
- Patil AP, Deka PC (2017) Performance evaluation of hybrid wavelet-ANN and wavelet-ANFIS models for estimating evapotranspiration in arid regions of India. *Neural Comput Appl* 28(2):275–285
- Rezakazemi M, Dashti A, Asghari M, Shirazian S (2017) H2-selective mixed matrix membranes modeling using ANFIS, PSO-ANFIS, GA-ANFIS. *Int J Hydrog Energy* 42(22):15211–15225
- Rumelhart DE, Hinton GE, Williams RJ (1986) Learning internal representation by error propagation. In: Rumelhart DE, McClelland JL (eds) *Parallel distributed processing*, vol 1. MIT, Cambridge
- Sayed T, Tavakolie A, Razavi A (2003) Comparison of adaptive network based fuzzy inference systems and b-spline neuro-fuzzy mode choice models. *J Comput Civ Eng ASCE* 17(2):123–130
- Shamshirband S, Amirmojahedi M, Gocić M, Akib S, Petković D, Piri J, Trajkovic S (2016) Estimation of reference evapotranspiration using neural networks and cuckoo search algorithm. *J Irrig Drain Eng* 142(2):04015044
- Tabari H, Kisi O, Ezani A, Talaei PH (2012) SVM, ANFIS, regression and climate-based models for reference evapotranspiration modeling using limited climatic data in a semi-arid highland environment. *J Hydrol* 444:78–89
- Takagi T, Sugeno M (1985) Fuzzy identification of systems and its applications to modeling and control. *IEEE Trans Syst Man Cybern* 15:116–132



- Termeh SVR, Kornejady A, Pourghasemi HR, Keesstra S (2018) Flood susceptibility mapping using novel ensembles of adaptive neuro fuzzy inference system and metaheuristic algorithms. *Sci Total Environ* 615:438–451
- Tsukamoto Y (1979) An approach to reasoning method. In: Gupta M, Ragade RK, Yager RR (eds) *Advances in fuzzy set theory and applications*. Elsevier, Amsterdam, pp 137–149
- Walls S, Binns AD, Levison J, MacRitchie S (2020) Prediction of actual evapotranspiration by artificial neural network models using data from a Bowen ratio energy balance station. *Neural Comput Appl*. <https://doi.org/10.1007/s00521-020-04800-2>
- Wen X, Si J, He Z, Wu J, Shao H, Yu H (2015) Support-vector-machine-based models for modeling daily reference evapotranspiration with limited climatic data in extreme arid regions. *Water Resour Manag* 29(9):3195–3209
- Wu L, Huang G, Fan J, Ma X, Zhou H, Zeng W (2020) Hybrid extreme learning machine with meta-heuristic algorithms for monthly pan evaporation prediction. *Comput Electron Agric* 168:105115
- Yuan X, Chen C, Lei X, Yuan Y, Adnan RM (2018) Monthly runoff forecasting based on LSTM–ALO model. *Stoch Environ Res Risk Assess* 32(8):2199–2212
- Zhu S, Heddad S et al (2019) Modeling daily water temperature for rivers: comparison between adaptive neuro-fuzzy inference systems and artificial neural networks models. *Environ Sci Pollut Res* 26(1):402–420



# Mapping shoreline change using machine learning: a case study from the eastern Indian coast

Lalit Kumar<sup>1</sup> · Mohammad Saud Afzal<sup>1</sup> · Mohammad Mashhood Afzal<sup>2</sup>

Received: 1 April 2020 / Accepted: 11 June 2020 / Published online: 30 June 2020  
© Institute of Geophysics, Polish Academy of Sciences & Polish Academy of Sciences 2020

## Abstract

The continuous shift of shoreline boundaries due to natural or anthropogenic events has created the necessity to monitor the shoreline boundaries regularly. This study investigates the perspective of implementing artificial intelligence techniques to model and predict the realignment in shoreline along the eastern Indian coast of Orissa (now called Odisha). The modeling consists of analyzing the satellite images and corresponding reanalysis data of the coastline. The satellite images (Landsat imagery) of the Orissa coastline were analyzed using edge detection filters, mainly Sobel and Canny. Sobel and Canny filters use edge detection techniques to extract essential information from satellite images. Edge detection reduces the volume of data and filters out worthless information while securing significant structural features of satellite images. The image differencing technique is used to determine the shoreline shift from GIS images (Landsat imagery). The shoreline shift dataset obtained from the GIS image is used together with the metrological dataset extracted from Modern-Era Retrospective analysis for Research and Applications, Version 2, and tide and wave parameter obtained from the European Centre for Medium-Range Weather Forecast for the period 1985–2015, as input parameter in machine learning (ML) algorithms to predict the shoreline shift. Artificial neural network (ANN), k-nearest neighbors (KNN), and support vector machine (SVM) algorithm are used as a ML model in the present study. The ML model contains weights that are multiplied with relevant inputs/features to obtain a better prediction. The analysis shows wind speed and wave height are the most prominent features in shoreline shift prediction. The model's performance was compared, and the observed result suggests that the ANN model outperforms the KNN and SVM model with an accuracy of 86.2%.

**Keywords** Shoreline change · Image processing · Artificial neural network · Edge detection · Machine learning

## Introduction

The shoreline change is a complex phenomenon that occurs due to the dynamic interaction of the ocean with the ground surface. The shoreline shifts are subjected to hydrodynamic forces observed in the sea that moves the sand, namely currents, winds, and waves. The interaction between landform and ocean causes erosion and accretion environment that exhibit in variable patterns (Morton 1996). The shoreline erosion is subjected to sea level rise (Bruun 1962). It relies

on the rate at which sediment is deposited and eroded from the seashore (Esteves et al. 2006). Shoreline also changes due to seasonal variation. The accretion of shoreline is gradual due to low waves of energy during the summer season. In contrast, rapid shoreline erosion is observed due to high storm waves during the winter season. However, seawalls and other stable structures limit the natural coastal cycle shifts up to a small extent, subjecting comparable losses to floating coastal areas. Also, the beaches, dunes, salt marshes, and estuaries can be in danger without sediment transport due to the introduction of artificial coastal structures.

Sediment transport is an important feature of nature which is affected by climate, magnitude, and direction of wave, wind energy, and many other factors. The sediment size is an essential factor for the movement of sediments through the river stream (Gazi et al. 2019; Gazi and Afzal 2020; Afzal et al. 2020). Therefore, sediments are classified based on their mode of transportation in the coastal region.

✉ Mohammad Saud Afzal  
saud@civil.iitkgp.ac.in

<sup>1</sup> Department of Civil Engineering, Indian Institute of Technology, Kharagpur, West Bengal, India

<sup>2</sup> Birlasoft Limited, Tower 3 Assotech Business Cresterra Plot no 22, Noida-Greater Noida Expy, Sector 135, Noida, Uttar Pradesh, India

The sediments are also deposited to the shore of the ocean by the river. Further, the transportation of sediment exists in the direction of onshore, offshore, and longshore. Sediment transport influences the rate of deposition and erosion, which results in the shoreline shift. The shift in shoreline directly affects non-living as well as the living matter and poses social, economic, and environmental threats (Small and Nicholls 2003; Dada et al. 2019). In one of the studies, Kumar et al. (2010) illustrated Coastal Vulnerability Index (CVI) for Orissa (now called Odisha) by using eight parameters. They demonstrated the utility of remote sensing data for coastal risk analyses, in situ measurements, computational simulation, and GIS research software. Thus, continuous monitoring and mapping of shoreline shifts are required to avoid the above consequences.

The shift in shoreline and its future prediction is essential in coastal and marine transportation, coastal zone management (CZM), and sediment transport. Initially, ground surveying techniques were used to study the coastline shift from 1807 to 1927. Later, the aerial photography was being introduced and used till 1980. The aerial photographic technique has some limitations, like transferring the information from images to maps and more time-consuming. Also, the black and white images (non-digital) of aerial photography were difficult to understand and interpret (De Jong and Van der Meer 2007). Hence, the photographic image data were replaced by the Landsat imagery and other remote sensing digital data since 1972. The digital data and new image processing techniques have eased the task of the researchers to map the shoreline shift (Alesheikh et al. 2007). However, historical time-scale datasets along with Geographical information system (GIS) imagery were also used in the prediction of shoreline erosion (Bagheri et al. 2019).

The remote sensing and GIS data have played a significant role in analyzing the shoreline shift (Howarth and Wickware 1981; White and El Asmar 1999; Chalabi et al. 2006; Zhang and Wang 2010). The basic methodology was to incorporate the physics between satellite images (multiple years) and observation sites data like temperature and pressure. The combination of satellite data and measurement data from the observation center was used in different models like linear regression (LR), end point rate (EPR), and minimum description length (MDL) to detect change in shoreline (Dolan et al. 1991; Li et al. 2001; Mukhopadhyay et al. 2012; Nandi et al. 2016). The EPR model is based on finding the ratio of change in shoreline distance with time elapsed. In contrast, the LR method is based on establishing the linear relation between the variables to predict the trends. These models predicted the shift by developing the mathematical model that specifies the interaction among the parameters assuming a specific distribution. The relationship between the parameters may have linear or nonlinear depending on the physics-based data generated at the

recording observation station (Gunawardena et al. 2009; Larson et al. 2000). Therefore, the parameters contributing to the coastal retreat require further exploration (Ramesh et al. 2017).

Previous studies showed that remote sensing data could be used in environmental monitoring programs to observe surface change phenomena over time. Further, the remotely sensed (aerial photography) images and field survey data were incorporated to observe beach width change in Gonghyunjin and Songjiho Beaches, South Korea (Kim et al. 2013). They also demonstrated the impact of artificial structures on shoreline shift. Further, the airborne Lidar and Unmanned aerial vehicles (UAVs) tool were used to observe remote sensing data and predict topographical changes in Uljington, Korea (Lee et al. 2019). The shoreline change extracted was statistically quantified using net shoreline movement (NSM) and linear regression rate (LRR) in the digital shoreline analysis system (DSAS). The Lidar and UAVs, generated digital surface maps (DSM), were compared to quantify morphological change. They reveal Lidar and UAVs-based digital surface maps can be used for coastal zone management.

Remote sensing is a geographic analysis tool capable of producing large quantities of data in the spatial, temporal, and spectral domains (Dellepiane et al. 2004; Alexakis et al. 2012; Nowakowski 2015). Hence, artificial intelligence (AI) techniques were used with remote sensing to interpret the large quantities of data for the image analysis process. The AI techniques are capable of modeling the mathematical models that have inherited non-linearity. The researchers have used different AI algorithms to analyze the shoreline shift (Ahmadian and Simons 2018; Hashemi et al. 2010). Ahangarha et al. (2019) proposed a procedure to determine land surface changes within the semiarid wetland and surrounding upland areas. They used a combination of hyperspectral images and machine learning (ML) algorithms to detect the shift in the land surface. Similarly, Kesikoğlu et al. (2020) observed the seasonal coastline shifts from Satellite images (Landsat 8) data using machine learning (ANN and KNN) model. They introduced a spatial pixel-based and object-based image classification technique to recognize changing areas at the coastline. Further, Harley et al. (2019) proposed a low-cost method to monitor shoreline change using smartphone images collected from social media platforms. They used an edge detection technique to find shoreline change on geo-rectified images. Their proposed approach provides almost similar prediction results of shoreline change as compared to GIS satellite images.

Ryan et al. (1991) used neural networks tool to delineate the shorelines. They used neural network with back propagation algorithm to categorize small blocks of image data. The image processing techniques were used in order to delineate the shoreline down to the pixel level. Further, Tsekouras

et al. (2018) investigated the shoreline realignment along an urban beach using a novel Hermite polynomial neural network model. They used reef morphology and wave forcing parameters for the modeling and obtained shoreline position with high spatiotemporal resolution. Similarly, Peponi et al. (2019) used a combination of ANN and GIS images to predict the coastal erosion in Costa da Caparica, Lisbon, Portugal. The impact of coastal erosion and shoreline change was obtained using the GIS-ANN model with better accuracy. Further, an advanced ML algorithm, DeepUNet (deep learning algorithm), was modeled to detect the shoreline change (Dickens and Armstrong 2019). The prediction of seafloor depth was observed using recurrent neural network (RNN) model at Mariana Islands, Marshall Islands, Guam, and Wake Island (USA). Their results did not match with the International Hydrographic Organization (IHO) standards as compared to the interpolated nautical chart data.

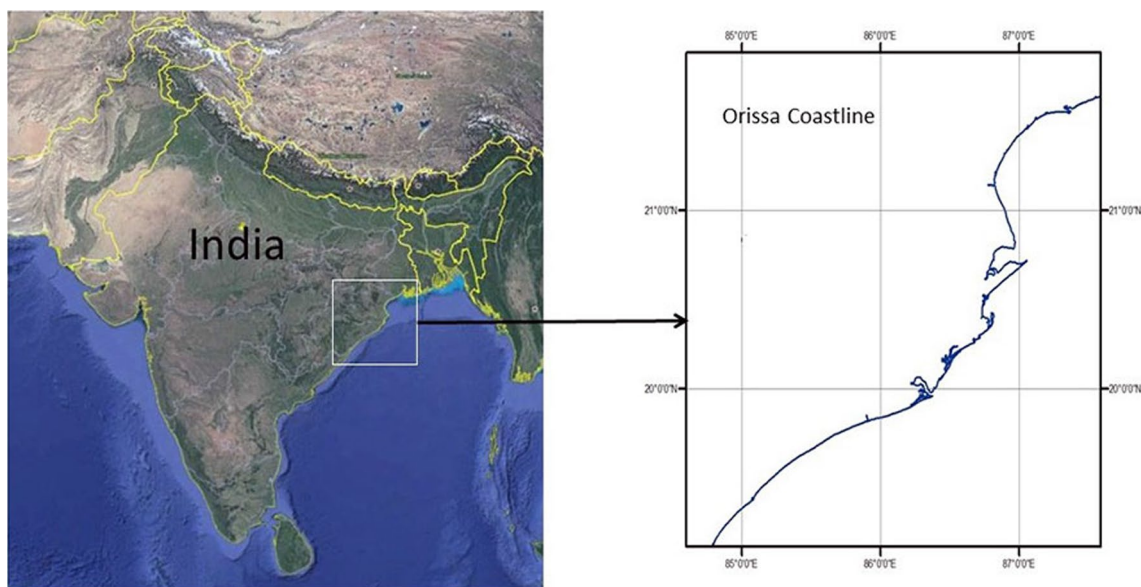
The machine learning (ML) tool had widely been used to predict various aspects of coastal and civil engineering in recent decades. Pierini et al. (2013) compared the machine learning (ANN) model prediction performance of hourly tidal level variations with a numerical model (MOHID model) at Puerto Belgrano (Argentina). The ML model outperforms the numerical model performance with better accuracy. Recently, Khaledian et al. (2020) used the machine learning (ANN and SVM) model to estimate the water level of the Caspian Sea. Montaña et al. (2020) compared the performance of 19 different numerical models with the ML model at Tairua Beach, New Zealand. The ML model shows accurate forecasting of shoreline change.

The conventional methods used for the study of shoreline were more complex, time-consuming, and require

more human resources through ground survey methods. The advancements in the GIS have made it feasible to overcome the shortcomings of the traditional survey methods. The use of ML techniques in a variety of coastal problems has rapidly increased over the past few years since ML algorithms can be highly effective predictors, can be used as part of larger models, and can provide physical insight. The present study focuses on both short-term and long-term shoreline changes using ML techniques to predict the shoreline shift of the Orissa coastline by utilizing satellite and physics-based data. This study sheds further insights into implementation of ML techniques to predict the realignment of coastline in general and extends present state of the art of coastline change prediction which had been previously been done using conventional methods. To the best of the authors knowledge, the present study is the first of its kind that uses ML models to predict the shoreline shifts in the entire Orissa coastline.

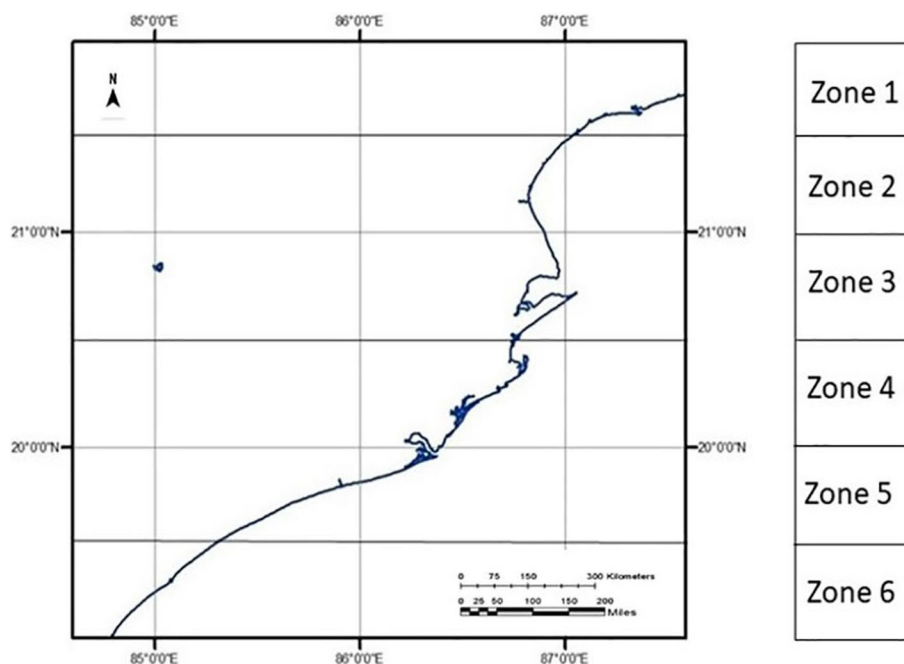
## Study region

The Orissa coastline is used as the study region to analyze the shoreline change. The Orissa is located in eastern India as shown in Fig. 1. The Orissa has 485 kilometers (301 miles) long coastline along the Bay of Bengal on its east, from Balasore to Ganjam. The coastline was divided into six zones similarly as Ramesh et al. (2017) as presented in Fig. 2. The region experiences four seasons: winter, pre-monsoon, southwest monsoon, and northeast monsoon. The Orissa coast is disaster-prone mainly due to flooding, surges, coastal inundation tropical cyclone, and tsunamis.



**Fig. 1** Study region: Orissa coastline

**Fig. 2** Zone-wise division of Orissa coastline



The growing appeal for trade and overseas investments has accelerated environmental changes in the state which has worsened in recent times due to the fluctuating weather conditions (heatwaves, cyclones, droughts, and floods). As a result, Orissa coast experiences severe erosion and deposition at different locations, which indicates the need for coastal management. Monalisha and Panda (2018) analyzed the eastern Indian coast Ganjam from 1972 to 2016 and revealed the shoreline had experienced a shift and accretion of approximately 5 km<sup>2</sup> and 1.6 km<sup>2</sup>, respectively. Similarly, Ramesh et al. (2017) reported that many villages of the Mahanadi deltaic coast were evacuated due to high erosion rates. Hence, it is necessary to understand the accretion and erosion rates over the coming years in this region.

## Research data

The Orissa coastline satellite images are obtained from United States Geological Survey (USGS) (USGS 2017) data for the duration between January 1985 and December 2015 on a monthly basis. The physics-based dataset is collected from the Modern-Era Retrospective analysis for Research and Applications, Version 2 (MERRA-2) (MERRA-2 2017; Kennedy et al. 2011; Gelaro et al. 2014; Bosilovich et al. 2017; Valipour and Tian 2018; Shen et al. 2019, 2020), and the European Centre for Medium-Range Weather Forecast (ECMWF) (ECMWF 2018; Hsu and Hoskins 1989; Bazile et al. 2017; Wang et al. 2019; de Rosnay et al. 2020). The comparison of the specifications of MERRA-2 and ERA-Interim (ECMWF) reanalysis data is presented in Table 1

MERRA-2 (Gelaro et al. 2017) is the 2nd version of the Modern-Era Retrospective analysis for Research and

**Table 1** A comparison of the reanalysis data

Name	ERA Interim	MERRA-2
Source	ECMWF	NASA GMAO
Time range	1979–present	1980–present
Assimilation	4D-VAR	3D-VAR with incremental update
Model resolution	TL255L60 and N128 reduced Gaussian	Native cube sphere grid output is interpolated to 5/8 lon., 1/2 lat.- deg.; 72 sigma levels
Dataset resolution	User defined, down to 0.75° × 0.75°	5/8 lon., 1/2 lat. degree, 42 pressure levels down to 0.67° × 0.5°
Dataset observed	Wave height, wave period, wave direction, swell height, swell direction, and tide	Temperature, relative humidity, pressure, wind speed, and wind direction

Applications, Version 2 (MERRA) developed by NASA's GMAO. MERRA-2 is generated to supplement the previous MERRA reanalysis data and addresses the later shortcomings in the assimilation of the newest satellite data sources (Rienecker et al. 2011). The MERRA-2 (Reichle et al. 2011) reanalysis data have been developed with a spatial resolution of  $0.67^\circ \times 0.5^\circ$  (Bosilovich et al. 2008) at NASA's Goddard Space Flight Centre. The latest version maintains some of its predecessor's core features, such as the spatial and temporal resolutions and the 3D Var 6-h update cycle. The parameters selected from the MERRA-2 are temperature (K), relative humidity (percent), pressure (bar), wind speed (m/s), and wind direction.

ERA-Interim is the 4th-generation reanalysis dataset of ECMWF, which follows the ERA 15 and ERA 40 dataset. ERA-Interim uses a 12-h, 4-dimensional variance analysis (4D Var) focused on the ECMWF with an efficient estimate of differences in satellite radiance results (Var BC). It enhances the correction of satellite observations (Dee et al. 2011). The parameters selected from ERA-Interim (ECMWF) are tides and wave parameters such as wave height, wave period, wave direction, swell height, and swell direction. The dataset contains a spatial resolution of  $0.75^\circ \times 0.75^\circ$ .

The tidal range in the study region varies from 0.7 m during neaps to 2.8 m during springs (The Indian Tide Tables-Part 1,1995: Indian and Selected Foreign Ports 1994). In Odisha coast, the mean significant wave height ranges between 1.25 and 1.40 m for the wave period of 6–9 s, mostly plunging from June to December and surging from January to May. The mean significant swell height is 1.33 m in the eastern Indian coast. The Orissa coastline poses the maximum temperature ranges between 35 and 40° C in summer and the low temperatures 3–4° C in winter. The relative humidity varies between 33 and 85%, whereas the pressure varies between the ranges of 998 and 1016 bar. The wind speed in the Orissa coastline varies in between 3 and 35 km/h throughout the year (with an average speed of 18 km/h). The wind is most often generated from the south with a peak percentage of 87%.

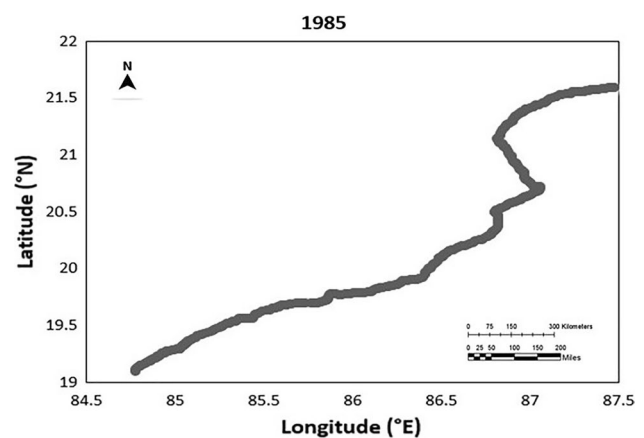
The sediment transport also plays a vital role in the detection of shoreline change of Orissa as the Mahanadi River at the delta bears an annual gross sediment load of 29.77 million tons per year (Ramesh et al. 2017). A considerable amount of sediment transport occurs in the coastline of Orissa. But due to the unavailability of sediment transport data, the sediment transport parameter is not considered in the present study. The choice of the reference line (shoreline) is the most critical part of shoreline change detection. The dune baseline is preferred over high water line due to its direct relationship with the tide and wave in the present study region (Elko et al. 2002; Houser et al. 2008; Stockdon et al. 2009; Suarez et al. 2012).

## Methods and numerical models

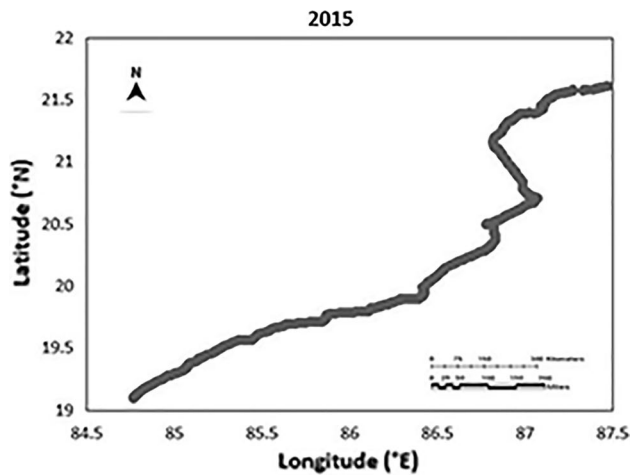
The significant local change of the intensity occurs at the boundary of two different regions in an image called the edges. First, the images were smoothed, and the noises were suppressed without destroying the true edges. In the second step, the sharpening filter was used to improve the quality of edges in the image. The third step was to detect the edge pixel, which should be either retained or discarded as noise depending on the threshold criterion for detection. Finally, localization determines the exact location of an edge. The coastline can be observed precisely by considering the first nonzero pixel from the right side as it represents the edge of the coastline. Moving down from one column to another, the first nonzero pixels from the right side represent the coastline. Figure 3 represents the Satellite image of Orissa coastline after edge detection in April 1985. However, Fig. 4 represents the Satellite image of Orissa coastline after edge detection in April 2015.

After superimposing the two satellite images (Landsat imagery), the coastline shift can be calculated by finding the number of black pixels between two consecutive white pixels in a row. Figure 5 represents the zone-wise shoreline shift (1985–2015) after superimposing the satellite images.

The satellite images of Orissa were analyzed using Sobel and Canny edge detection technique to observe shoreline shift. The Canny edge detection technique has previously been shown to give higher accuracy in detection of shoreline shift and less execution time compared with Sobel edge detection technique (Acharjya et al. 2012; Vijayarani and Vinupriya 2013). Similarly, the canny edge detection outperforms Sobel edge detection technique and gives better edge detection results with higher accuracy in



**Fig. 3** Satellite image of Orissa coastline after edge detections in April 1985

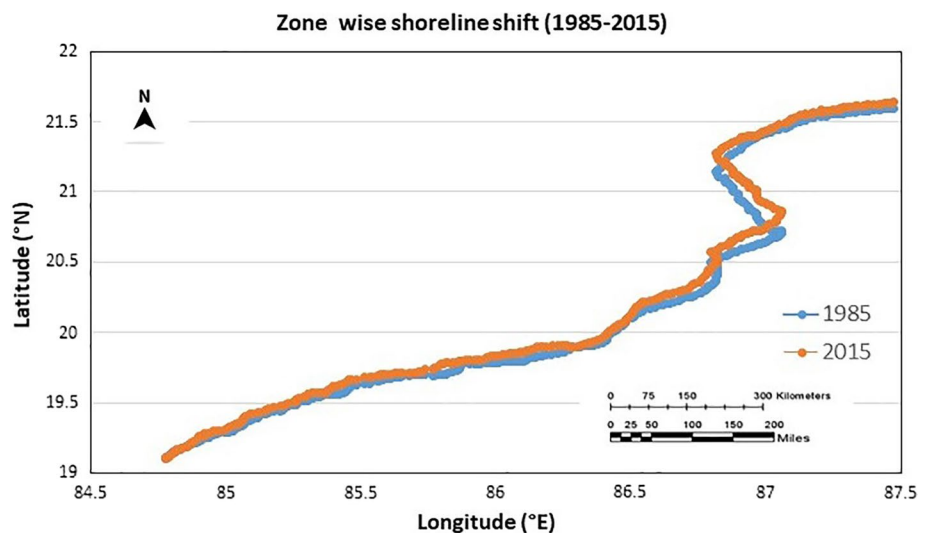


**Fig. 4** Satellite image of Orissa coastline after edge detections in April 2015

the present study. Therefore, the coastline shift obtained using the Canny edge detection techniques was used as an input in ML models.

The data provided by MERRA-2 and ERA-Interim (ECMWF) are reanalysis data. The reanalysis data are preprocessed data that do not require any cleansing or feature engineering. The dataset obtained from MERRA-2 and ECMWF is in NetCDF format. Further, the dataset is extracted from the NetCDF format for the Orissa coastline in excel file format using MATLAB 2019a. Thereafter, the data were normalized using a standard scale before being used as a model input parameters in order to avoid the problem of scaling. The normalization process transforms the data having a distribution that exhibits the mean value of zero and a standard deviation of one. The feature matrix was normalized using Eq. 1.

**Fig. 5** Satellite image of Orissa coastline after superimposing the image of 2015 over 1985



$$X_{\text{normalized}} = \frac{X - \mu}{\sigma} \quad (1)$$

where  $\mu$  = mean (X) and  $\sigma$  = standard deviation (X)

Further, the ML models were trained from tides, shore-line shifts, and the wave parameters from ERA-Interim (ECMWF). Subsequently, the trained ML models, viz. artificial neural network (ANN), K-nearest neighbor (KNN), and the support vector machine (SVM), are used to model and predict the future coastline shift. The k-fold cross-validation technique is used to estimate how the models are expected to perform when used to predict the test dataset. Figure 6 shows a schematic flowchart that illustrates the methodology used in the present study.

### Sobel edge detection

The Sobel filter was introduced by Sobel–Feldman (Sobel and Feldman 1968). It is a discrete differential operator, which computes an estimate of the feature intensity function gradient. The corresponding gradient vector at each point in the image is the product of the Sobel–Feldman operator. The Sobel–Feldman operator is generally used for converting the image in horizontal and vertical directions with a small, separable, and integer-valued filter and is therefore relatively inexpensive computationally. On the other hand, the gradient approximation generates relatively coarse images, especially with respect to high-frequency image variations.

Sobel edge detection technique detects the edges by seeking the maximum and minimum in the first derivative of the image (Vincent et al. 2009; Saluja et al. 2013). The Sobel filter performs two-dimensional (2D) gradient measurement on images returning the maximum gradient edge. For edge detection operations, pairs of horizontal

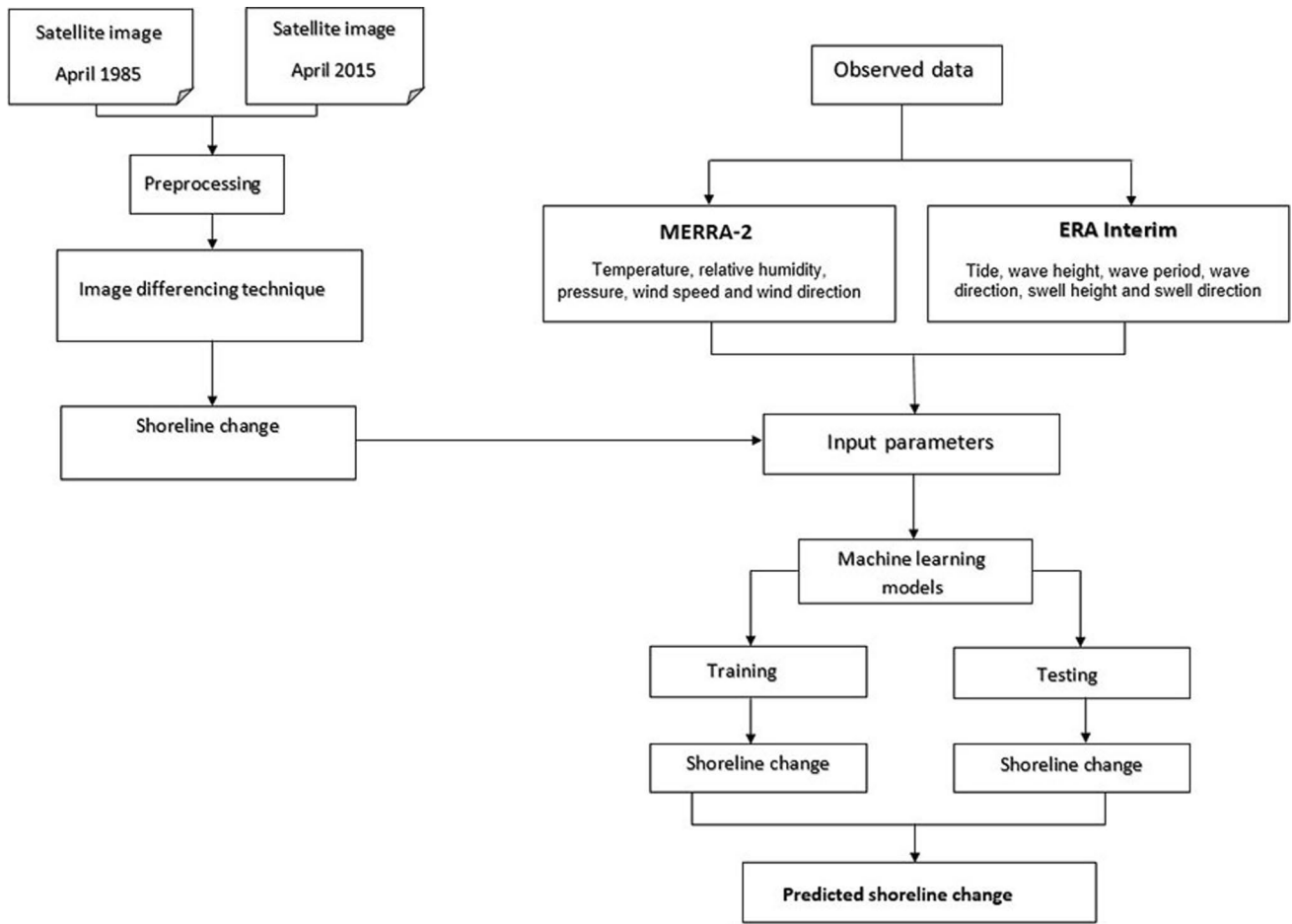


Fig. 6 Flowchart of the methodology

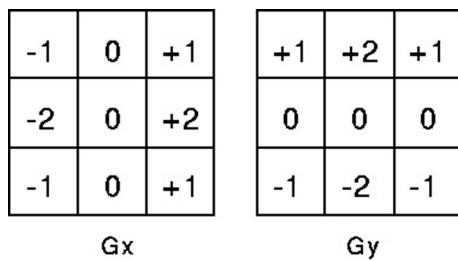


Fig. 7 Horizontal and vertical gradient matrices of Sobel filter

and vertical gradient matrices of dimensions  $3 \times 3$  are used and presented in Fig. 7.

The gradient is a vector, and their components determine the change of rapid pixel values with respect to distance in the x and y directions. The components of the gradient can be calculated using Eqs. 2 and 3:

$$\frac{\partial f(x, y)}{\partial x} = \Delta x = \frac{f(x + dx, y) - f(x, y)}{dx} \tag{2}$$

$$\frac{\partial f(x, y)}{\partial y} = \Delta y = \frac{f(x, y + dy) - f(x, y)}{dy} \tag{3}$$

where dx and dy show the distance measured along x and y directions, respectively. dx and dy can be considered in terms of numbers of the pixel between two points in discrete images. The pixel coordinate at a point (i, j) is presented in Eqs. 4 and 5, where dx = dy = 1 (pixel spacing).

$$\Delta x = f(i + 1, j) - f(i, j) \tag{4}$$

$$\Delta y = f(i, j + 1) - f(i, j) \tag{5}$$

To find the presence of discontinuity in gradient, the change in gradient at (i, j) is calculated by finding out the magnitude using Eq. 6.

$$M = ((\Delta x^2) + (\Delta y^2))^{\frac{1}{2}} \tag{6}$$



Direction  $\theta$  is calculated using Eq. 7.

$$\theta = \arctan(\Delta x / \Delta y) \quad (7)$$

### Canny edge detection

John F. Canny introduced the Canny edge filter in 1986 (Canny 1986). It is an operator for edge detection which uses a multistage algorithm to detect a broad range of edges in images. This technique is developed to extract structural information from different objects of vision and significantly minimize the volume of data to be processed. It has been spread across various machine vision platforms commonly. Canny has also identified fairly close criteria for the implementation of edge detection on different vision systems. Thus, Canny edge detection can be implemented in a wide range of situations to address image processing requirements. The Canny edge detection technique locates the edge at the edge center with a low error rate, and no false edges are created during its execution.

Canny edge detection techniques detects the edges with suppression of noise (Green 2002; Shrivakshan and Chandrasekar 2012). It smoothens the image using a Gaussian filter that reduces noise using Eqs. 8 and 9 (Murthy et al. 2009).

$$g(m, n) = G_{\sigma}(m, n) * f(m, n) \quad (8)$$

$$G_{\sigma} = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{m^2 + n^2}{2\sigma^2}\right) \quad (9)$$

where  $\sigma$  stands for the parameter of Gauss filter, and it controls the extend of smoothing image.

Using gradient operations, Eqs. 10 and 11 were obtained.

$$M(m, n) = \sqrt{g_m^2(m, n) + g_n^2(m, n)} \quad (10)$$

$$\theta(m, n) = \tan^{-1}[g_n(m, n) / g_m(m, n)] \quad (11)$$

where the threshold value of  $M$  is obtained by Eq. 12, if  $M(m, n) > T_0$ ; otherwise, it is 0.

$$M_T(m, n) = M(m, n) \quad (12)$$

where  $T$  is selected in such a way that all edge elements were considered and noise is also suppressed up to the maximum extent.

### K-nearest neighbors (KNN)

KNN is an AI algorithm which accumulates all available cases and classifies new cases based on the equivalent

measure (distance function) (Satapathy et al. 2012). KNN algorithm is an easy-to-implement supervised AI algorithm which is used to solve classification as well as regression problem (Altman 1992). KNN algorithm is developed to perform statistical estimation and pattern recognition (Garg et al. 2019). A case is classified by a majority vote of its neighbors, with the case being assigned to the class most common amongst its K-nearest neighbors measured by a distance function as given in Eqs. 13, 14, and 15.

### Distance functions

$$\text{Euclidean} \Rightarrow \sqrt{\sum_{i=1}^n (x_i - y_i)^2} \quad (13)$$

$$\text{Manhattan} \Rightarrow \sum_{i=1}^n |x_i - y_i| \quad (14)$$

$$\text{Minkowski} \Rightarrow \left( \sum_{i=1}^n |x_i - y_i|^q \right)^{\frac{1}{q}} \quad (15)$$

Minkowski distance function is the generalized form of Euclidean and Manhattan distance function. Therefore, Minkowski distance function results were used in the present study. In this technique, positive integer  $k$  and a new sample are specified. The positive integer  $k$  in the dataset which is closest to the new sample is selected to find out the most common class, and their classification is then given to the new sample. KNN algorithm predicts output by calculating the similarities between an input sample and each training instance.

### Support vector machines (SVM)

SVM is an AI algorithm developed in Russia in the 1960s (Vapnik 1963; Vapnik and AY 1965). SVM was developed on large scale at AT & T Bell Laboratories by Vapnik and co-workers (Cortes and Vapnik 1995). SVMs are structured learning models that are used for classification and regression analysis problems. The SVM training algorithm builds a model for a specified set of training examples, allocating a new example to one or another subclass, rendering it a binary non-probabilistic linear classifier (Cortes and Vapnik 1995; Hsu et al. 2003; Dutta et al. 2020). Also, the regression problem can be solved using the SVM model.

SVM algorithms employ a number of mathematical functions which are defined as the kernel. The kernel's function uses the input data and transforms it into the required form.. Basically, it returns the internal product between two

points in an appropriate feature space. Different SVM algorithms use different types of kernels functions such as linear, nonlinear, polynomial, radial basis function, etc. (Gunn et al. 1998). In recent years, kernel methods have gained considerable attention due to the support vector machine's popularity. In many applications, kernel functions provide a simple connection between linearity and non-linearity for algorithms, which can be expressed as dot products (Fadel et al. 2016). A nonlinear classifier is created by applying the kernel functions to create maximum-margin hyperplanes. The appropriate choice of kernel function is an important choice selection that affects the overall accuracy. Hence, the choice of kernel is one of the important aspects while implementing SVM. The polynomial, radial basis, and sigmoid kernel functions are used in the present study.

In SVM algorithm, the radial basis function kernel (RBF) is often used for classification analysis. For two samples  $x(i)$  and  $x(j)$ , the RBF kernel function can be represented as feature vectors in an input space (Chudzian 2011), as presented in Eq. 16:

$$K(x^i, x^j) = \phi(x^i)^T \phi(x^j)^T = e^{-\gamma \|x^i - x^j\|^2} \quad (16)$$

where  $\gamma = \frac{1}{2\sigma^2}$  and  $\phi(x)$  is infinite dimensional for this kernel.

The polynomial kernel (on degree- $n$  polynomials  $x$  and  $y$ ) is represented as feature vectors in an input space over polynomial of the original variable. For polynomials of degree- $n$ , the polynomial kernel is defined as in Eq. 17 (Gunn et al. 1998).

$$K(x, y) = (x^T y + c)^n \quad (17)$$

where  $c$  is greater or equal to zero ( $c = 0$  means homogeneous).

The sigmoid kernel originates from the neural networks, where artificial neurons often use the bipolar sigmoid function as an activation function. For two samples  $x(i)$  and  $x(j)$ , the RBF kernel function can be written as Eq. 18.

$$K(x, y) = \tanh(ax^T y + c) \quad (18)$$

where  $a$  and  $c$  represent the slope and intercept, respectively.

Radial basis function kernel, polynomial kernel, and sigmoid kernel have been used for the prediction of shoreline shift in the present study (comparison is shown in results section).

### Artificial neural network (ANN)

ANN is an AI algorithm that consists of simple elements called artificial neurons (Puskarczyk 2019; Bouguerra et al. 2019). The neural network model uses the same approach as our brain does (Gatys et al. 2015). It is a highly generalized

form of linear regression, resulting in extremely complex and nonlinear interactions between the input data and the output (Piasecki et al. 2018). The neuron receives input data from the source, changes their activation, and produces output as shown in Fig. 8.

The obtained output is completely dependent on the input data and their activation (Schalkoff 1997). Each neuron in the hidden layer gets weighted inputs (output of the previous layer) plus bias from each neuron in the previous layer, as presented in Eq. 19.

$$Z_i = \left( \sum_{k=1}^{N_j-1} X_k^{j-1} * W_{k,i} - b_k \right) \quad (19)$$

where  $Z$  = output,  $X$  = input parameter,  $W$  = weight, and  $b$  = bias

In this study, the activation function  $\tan(h)$  in exponential form is represented as in Eq. 20:

$$f(Z_i) = \frac{e^{(Z_i)} - e^{(-Z_i)}}{e^{(Z_i)} + e^{(-Z_i)}} \quad (20)$$

The summation is passed along the activation function to generate the output of the node, which is calculated as  $Y_i = f(Z_i)$

ANN involves three important steps: training, validation, and testing. In the training step, the network is trained by adjusting the weights. The second step validation is necessary to avoid over-fitting of the model. The validation set is directly not used for weights adjustment, but it is used to find out the optimum number of hidden layers and also decide the termination point. The third step is testing, which is used to check the prediction ability of the network.

The order or arrangements of node in a layer of neural network is defined by neural network architecture (Arce-Medina and Paz-Paredes 2009). The feed-forward neural network is used in the present study, which consists of an input layer, a series of hidden layers, and an output layer, each having various numbers of nodes (Valipour et al. 2012, 2013). The total number of nodes in layers is determined by the nature of a problem or its complication under

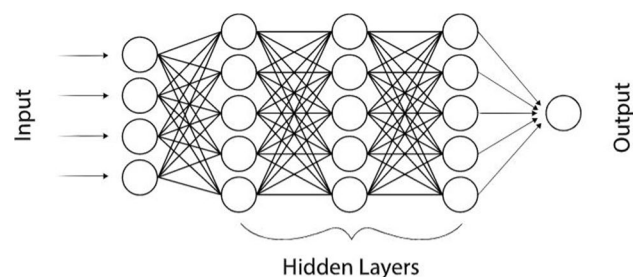


Fig. 8 A generalized schematic of a neural network

consideration. The weights are calculated when the dataset is trained in the neural network. Finally, the output is predicted from the input node, which acts as a distribution node.

One of the accepted ways to validate the ML model is using the train–test split method. In this method, the accuracy of the model is evaluated in terms of coefficient of determination ( $r^2$ ). The coefficient of determination ( $r^2$ ) is a statistical measure of how well the predictions approximate the real data points. The value of ( $r^2$ ) varies in between 0 and 1. A higher ( $r^2$ ) value indicates better prediction accuracy. In the present study, the dataset obtained from MERRA-2, ERA-Interim (ECMWF), and shoreline shift is divided into two parts based on the train–test split method. The first part contains 70% of the dataset, which is used for model training purposes, and the remaining 30% of the dataset is used as testing purposes. The training dataset is used to develop the ML model. However, the trained ML model is then tested with inputs from remaining 30% dataset to predict the shoreline shift of the testing dataset.

## Results and discussions

The coastline retreat was calculated by superimposing the two satellite images of the Orissa (India) coastline using edge detection filters. Image differencing technique was used to determine the coastline shift of Orissa (India). The shift is estimated by finding the number of black pixels between two consecutive white pixels in a row where a unit pixel corresponds to 200 meters. A significant shift of the Orissa coastline was observed in April 2015 (Landsat imagery) with respect to April 1986 (Landsat imagery). Subsequently, KNN, ANN, and SVM techniques are used to model the shift of Orissa (India) coastline in conjunction with other input parameters from MERRA-2 and ERA-Interim (ECMWF).

The weight parameters for the KNN model are set to be uniform such that every point in the neighborhood is weighted equally. The power parameter ( $p$ ) for the Minkowski metric was selected as 2 (two). Therefore, Minkowski metric is equivalent to the standard Euclidean metric (Hu et al. 2016). The accuracies of KNN models having a different number of neighbor parameters are presented in Table 2. The training and testing accuracy was

**Table 2** KNN model

Networks	Train accuracy	Test accuracy	No. of neighbors( $k$ )
KNN	0.879	0.830	5
KNN	0.852	0.835	10
KNN	0.845	0.845	100

found to be 0.879 and 0.830 for  $K = 5$  and 0.852 and 0.835 for  $K = 10$ . The training and testing accuracy for  $K = 100$  was 0.845 and 0.845. Table 2 shows that better result was obtained with  $k = 100$  in the present study.

In SVM model, the penalty parameter ( $C$ ) was set to one, and the degree of polynomial kernel was three. The gamma was set in the algorithm automatically, which uses  $1 / \text{number of features}$ , and if gamma is set to a scale, then it uses  $1 / (\text{number of features} * X.\text{var})$  as the value of gamma. The polynomial kernel with gamma function provides the training and testing accuracy of 0.855 and 0.845, respectively. The accuracy achieved by the SVM models with different kernels is presented in Table 3. The corresponding training and testing accuracy with RBF function was 0.835 and 0.845, respectively, and the sigmoid kernel function has the training and testing accuracy of 0.83 and 0.83, respectively. Table 3 shows the best result was obtained with polynomial kernel function in the present study.

An ANN model of (100,100) neurons having two hidden layers using tan hyperbolic as the activation function showcased the highest test accuracy. The ANN model predicts the shoreline shift with  $r^2$  of 0.862 and an accuracy of 86.2% as presented in Table 4. However, the stochastic gradient descent (SGD) and limited memory Broyden–Fletcher–Goldfarb–Shanno (LBFGS) optimization algorithms (Jan et al. 2002) showed a test accuracy of 73 and 75%, respectively. The highest test accuracy was manifested by Adam optimizing algorithm as represented in Table 4. It can be concluded that an ANN model of (100,100,2) neuron

**Table 3** SVM model

Kernel	Gamma function	Train accuracy	Test accuracy
POLY	Scale	0.855	0.845
RBF	Auto	0.835	0.845
Sigmoid	Auto	0.830	0.830

**Table 4** ANN model with Adam optimizer

Train accuracy	Test accuracy	Neurons	Function
0.845	0.845	12,12,2	Identity
0.846	0.847	12,12,2	Tanh
0.847	0.847	12,12,2	Relu
0.847	0.847	12,12,2	Logistic
0.846	0.847	12,2,5	Relu
0.847	0.847	50,50,5	Relu
0.845	0.845	50,50,2	Relu
0.857	0.857	50,50,2	Tanh
0.859	0.858	100,100,30	Tanh
0.863	0.862	100,100,2	Tanh
0.858	0.856	100,100,30	Relu

**Table 5** Comparison of accuracy scores of ANN, KNN, and SVM models

Model	Train accuracy	Test accuracy
KNN	0.845	0.845
SVM	0.855	0.845
ANN	0.863	0.862

**Table 6** Results obtained from ANN model

Year	Average yearly shift	Average shift in meter
1985	0	0
1986	0	0
1987	0	0
1988	0.008	1.661
1989	-0.009	-1.898
1990	-0.032	-6.487
1991	0.355	71.044
1992	-0.214	-42.827
1993	0.268	53.691
1994	-0.295	-59.098
1995	-0.378	-75.685
1996	-0.204	-40.772
1997	0.355	71.044
1998	-0.254	-50.448

is sufficient and accurate, while further increasing the number of neurons and layers, no significant change is observed in the test accuracy.

The comparison of results from the KNN, SVM, and ANN model is presented in Table 5. The result suggests that ANN model using Adam optimizing algorithm has

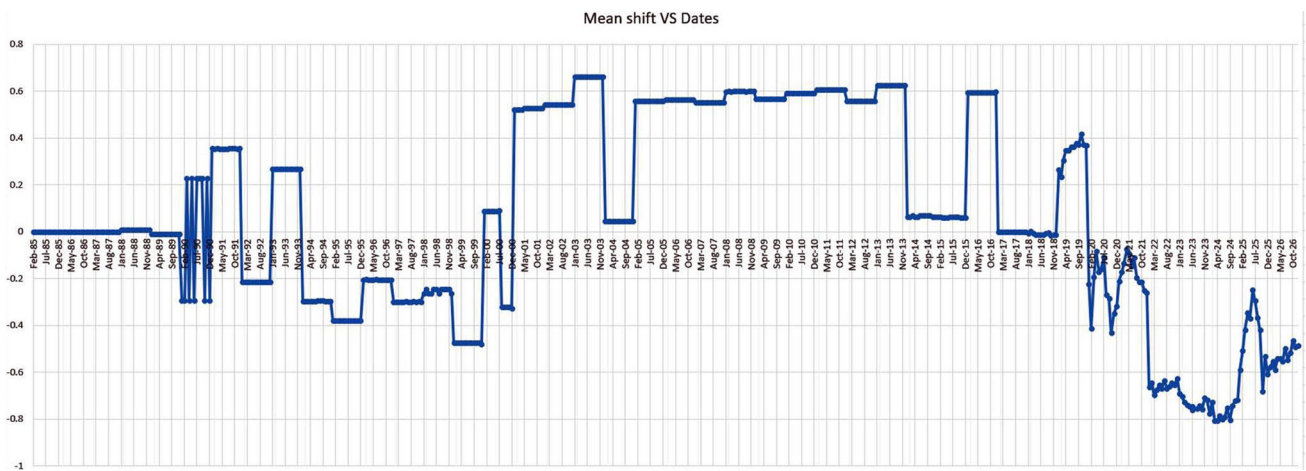
better accuracy compared to other models used in the present study. The ANN model outperforms the KNN and SVM model. It is evident that the monthly forecasting is improved significantly at every point location using the ANN model. The model can be further used in different research areas.

After applying the ANN model, the shifts of the shoreline for respective years were calculated for various points along the coastline. Subsequently, a mean shift was calculated as the total shift of the coastline from a reference time period of 1985, as presented in Table 6. It contains the average yearly shift and average shift in meter from 1985 to 1998. The mean shift in pixels is plotted against the time period, as presented in Fig. 9. The single unit pixel shift of the image is equivalent to a shift of 200 meters of shoreline.

The mean shift in future years (2018–2100) is calculated using a best fitting curve (yearly mean shift vs. time period) using Curve Pro-Expert Professional software as presented in Table 7 and shown in Fig. 10. A significant increase in the average shift in meter is observed in the prediction of shoreline shift in future years, as shown in Table 7 and Fig. 10. It is evident from the observations that the coastline would continue to erode at an increasing rate in future years.

The statistics of the curve are presented in Table 8. The coefficient of determination ( $r^2$ ) of the curve and correlation coefficient ( $r$ ) are 0.560 and 0.748, respectively. The correlation coefficient ( $r$ ) of the best fitting curve is closer to 1. This means it will give a better prediction of shoreline shift in future years.

The result of the ANN model to predict the zone-wise shift is presented in Table 9. The coastline was divided into six zones similar to Ramesh et al. (2017) as presented in Fig. 2. The ANN model result suggests that Zone 1 and Zone 2 will continue to accrete over the coming years at an increasing rate. However, Zone 3 would accrete at a decreasing rate. Zone 4 will continue to erode at an



**Fig. 9** Mean shift versus time period

**Table 7** Mean shift in meters

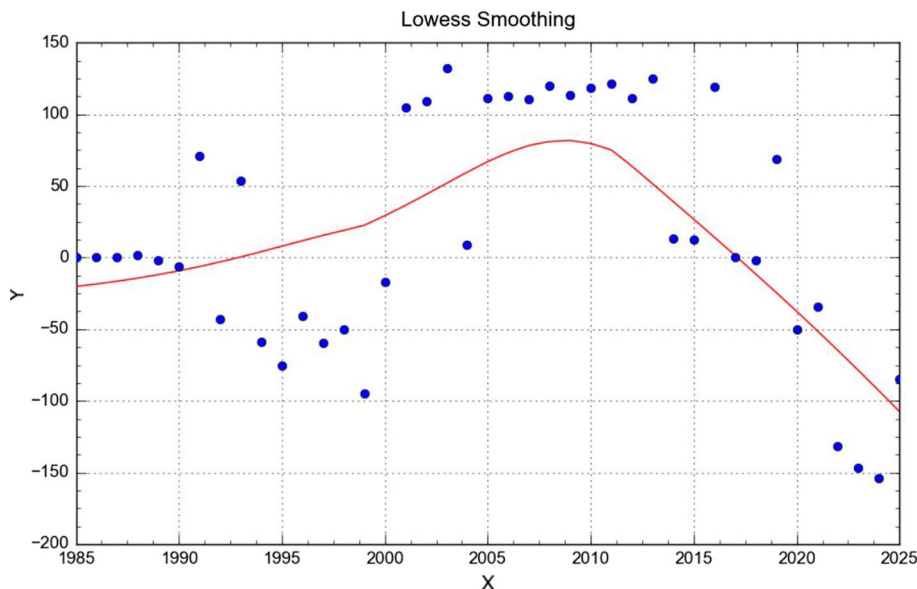
Year	Average shift in meter
2018	- 43.24
2025	- 107.52
2050	- 179.06
2075	- 250.67
2100	- 322.14

indicate that the wind speed and the wave height are the most prominent features in shoreline shifts.

**Comparison with the published literature and validation**

There is no such literature that covers the entire coastline of Orissa from 1975 to 2015 using the ML approach as used in the present study. In one of the published reports, Ramesh et al. (2017) analyzed the shoreline shifts of the Orissa coast-

**Fig. 10** Best fitting curve for mean shift versus time period



**Table 8** Best fitting curve statistics

Standard error	Correlation coefficient ( <i>r</i> )	Coefficient of determination ( <i>r</i> <sup>2</sup> )	Score
56.812	0.748	0.560	567

line from 1972 to 2010, where they used remote sensing techniques for their analysis. They have also not provided the quantitative results of the shoreline shift. However, Mishra et al. (2019) qualitatively validated their shoreline status at the Puri District of Orissa (1990–2015) against Ramesh et al. (2017). In a similar manner, the zone-wise erosion

**Table 9** Zone-wise shift in meters

Year	Zone 1	Zone 2	Zone 3	Zone 4	Zone 5	Zone 6
2016	34	39.62	113.46	- 798	18.86	3.60
2025	36.58	44	122	- 846	14	1.92
2050	41	80	94	- 894	- 3.8	- 2.8
2075	44	116	66	- 944	- 21.8	- 7.6
2100	49.96	152	38	- 994	- 40	- 12.6

increasing rate. However, Zone 5 and Zone 6 would erode at an increasing rate. The weights associated with the features of all the three ML models ( KNN, SVM, and ANN)

and accretion results of present study are validated with the National assessment report of Ramesh et al. (2017) as shown in Table 10. The results show that the results obtained in

**Table 10** Comparison of the present study of observed shoreline status against National assessment report by Ramesh et al. (2011) along the eastern Indian coast, Odisha

Zone	Length (km)	Present study (1985–2015)	National assessment report by Ramesh et al.(2011) (1972–2010)
Zone 1	87.96	Accretion	Accretion
Zone 2	52.61	Accretion	Accretion
Zone 3	83.55	Accretion	Accretion
Zone 4	58.95	Erosion	Erosion
Zone 5	136.48	Erosion	Erosion
Zone 6	60.85	Accretion	Accretion

the present study using ML technique match with that of Ramesh et al. (2017) qualitatively who used remote sensing techniques for their analysis.

The shoreline change assessment based on the ML technique reveals considerable temporal variability in the positions of shorelines along the eastern Indian coast of Orissa over the past 30 years (1985–2015). The shoreline shift observation documented in the present study is consistent mainly with recent shoreline change assessments in the eastern Indian coast, Odisha (Ramesh et al. 2017; Barman et al. 2014; Rajawat et al. 2015; Jangir et al. 2016). For instance, Ramesh et al. (2017) reported that out of total 480.4 km coastal length of Orissa, 46.8% is accretion, 36.8% is eroded, 14.38% is stable, and 2.04% is artificial coast during 1972 (Survey of India toposheet) and 2010 (Landsat-5 TM) that means accretion is dominant throughout these years similar to the present study. Further, Barman et al. (2014) reported that the shoreline of the Balasore District of Orissa varied from 1.4 to 3.75 m while in NDVI, the shift varied from 2.0 to 9.31 m due to accretion (1975–2013). They have also predicted the shoreline shift till 2030 with the linear regression technique and reported the accretion would take place at a higher rate in future years. Similarly, Zone 1 indicating Balasore poses similar future results and suggests that the accretion will take place in future years with an increasing rate. In one of the studies, Mukhopadhyay et al. (2012) used Landsat MSS and TM data in Puri (Odisha) during 1972–2010. They recorded high erosion rates in the northern part of Puri, near the Kushabhadra estuary and Chandrabhaga beach. In contrast, the southern portion of the shoreline between Chilika and Puri is substantially stable during 1972–2010. Similar results of shoreline shift were obtained by Jangir et al. (2016) as they documented that the high erosional trends (500–800 m) were observed in the northern part of Puri (Odisha) in comparison with the southern part (50–200 m) during 1972–2009 using the Survey of India toposheet (1972) and Landsat TM data (2009). Overall, the same shoreline shift trend has been observed in Zone 4 of this study, which poses high erosional trends. The significant

shoreline erosion results in Zone 4 are expected due to complex interactions between river flow, waves, and tides. However, Zone 5 undergoes erosion with a low erosion rate as compared to Zone 4. Similarly, Markose et al. (2016) analyzed the coastline of Ganjam (Orissa) during 1990–2014 using satellite imagery and reported that the 71.65% of the Ganjam coast undergoes accretion, whereas 28.35% coast falls under erosion. A similar trend is observed in Zone 6 of this study. There are no studies till date analyzed in Zone 2, Zone 3, and Zone 4 indicating Bhadrak, Kendrapara, and Jagatsinghpur districts, respectively. These zones can be considered separately and analyzed in future research works.

Machine learning (ML) algorithms are data-driven approaches which extract the interpretable information and knowledge from the available data resources. Thus, ML derives models that learn much more from the big data than the traditional data assimilation approaches can, while still respecting the evolving understanding of nature's law. The ML system's goal is to prepare a function that best maps inputs to outputs given the resources available. Therefore, it can be a valid assumption for using a certain number of input parameters as per the availability of the required dataset to ML models (Govindaraju 2000a, b). The human activities and coastal ecosystems relationships are difficult to understand since these activities generate multiple pressures acting simultaneously and often producing unexpected ecosystem responses (Halpern et al. 2008). In more recent times, human activities such as those related to land-use practices, the spread of urbanized areas (Valiela 2004) and the building of dams and offshore structure, have significantly reduced the delivery of fluvial sediments to the coastal systems and therefore altered the natural coastal processes of sedimentation (Coltori 1997; Gregory 2004; Simeoni and Corbau 2009; Ronco et al. 2010; Di Silvio and Nones 2014; Guerrero et al. 2015; Pescaroli et al. 2018; Varrani et al. 2019). The natural- and human-induced pressure is unexpected in the coastal region and presently not included in the current study, and therefore, the predicted value from ML approach may change from the actual results.

## Conclusion

The shoreline shift has been an important issue for the researcher in recent years due to continuous occurrence of both natural and anthropogenic events in the coastal region. The prediction of shoreline change is important with regard to coastal hazard assessment. The problem of coastal erosion in eastern Indian coast, Orissa (India), has increased due to high frequency and intensity of cyclones such as Helen (2013), Hudhud (2014), and Fani (2019), and repeated floods in recent years. The eastern Indian coastline, Orissa, is facing a high degree of erosion due to the

effect of natural phenomenon as well as human activities in that area. In the present study, ML algorithms (KNN, SVM, and ANN) were implemented to predict the realignment in shoreline along the eastern Indian coast of Orissa, India. A significant shoreline shift was observed using Canny filter (image differencing technique) in multiresolution satellite images (Landsat imagery) of 1985 and 2015. The reanalysis data from MERRA-2, ERA-Interim (ECMWF), and shoreline shift data were used to model ML algorithms. The ML algorithms were applied to calculate and forecast shoreline changes along the Orissa coastline. The ANN model outperforms KNN and SVM algorithm in shoreline prediction with an accuracy of 86.2%. The shoreline observed in the present study was dynamic and had uncertainties all along the coast. The study reveals the change in the shoreline of Orissa in future decades. The weights associated with the features of all the three ML models (KNN, SVM, and ANN) indicate that the wind speed and the wave height are the most prominent features in shoreline shifts.

Hence, it can be concluded that the precise and accurate prediction of the shoreline change can be observed using ML techniques. The above results state that the shoreline of the eastern coast, Orissa (India), will experience a reduction in the shoreline at a rate of 2.61 m/year while considering of the shoreline of 2018 as a reference along the coast.

The results of the study can be used as management tools for shorelines protection to avoid economic losses in the future. Vegetation, marshes, and stone groin can be provided to safeguard seashore in the coastline area. Sea level rise along the coast, cyclones, and repeated flood must be continuously investigated to manage or reduce the loss in future.

**Acknowledgements** We thank Mr. Sobhit and Mr. Satish Yadav (B. Tech students of IIT Kharagpur) for the assistance in writing code. This work was carried out as a part of the project titled “Predictive Tool for Arctic Coastal Hydrodynamics and Sediment Transport” funded by the National Centre for Polar and Ocean Research (NCPOR). Authors also acknowledge support by SRIC, IIT Kharagpur, under the ISIRD project titled “3D CFD Modeling of the Hydrodynamics and Local Scour Around Offshore Structures Under Combined Action of Current and Waves.”

**Funding** Funding was provided by Sponsored Research and Industrial Consultancy (Grant No. IIT/SRIC/CE/MOS/2017-18/200) and Ministry of Earth Sciences (Grant No. NCPOR/2019/PACER-POP/OS-02).

## Compliance with ethical standards

**Conflict of interest** The authors declare no conflicts of interest in the current paper.

## References

- Acharjya PP, Das R, Ghoshal D (2012) Study and comparison of different edge detectors for image segmentation. *Glob J Comput Sci Technol* 12:29–32
- Afzal MS, Bihs H, Kumar L (2020) Computational fluid dynamics modeling of abutment scour under steady current using the level set method. *Int J Sediment Res* 35:355–364
- Ahangarha M, Seydi ST, Shahhoseini R (2019) Hyperspectral change detection in wetland and water-body areas based on machine learning. In: International archives of the photogrammetry, remote sensing & spatial information sciences, geospatial conference 2019—joint conferences of SMPR and GI research, vol XLII-4/W18, pp 19–24
- Ahmadian AS, Simons RR (2018) Estimation of nearshore wave transmission for submerged breakwaters using a data-driven predictive model. *Neural Comput Appl* 29(10):705–719
- Alesheikh AA, Ghorbanali A, Nouri N (2007) Coastline change detection using remote sensing. *Int J Environ Sci Technol* 4(1):61–66
- Alexakis DD, Agapiou A, Hadjimitsis DG, Retalis A (2012) Optimizing statistical classification accuracy of satellite remotely sensed imagery for supporting fast flood hydrological analysis. *Acta Geophys* 60(3):959–984
- Altman NS (1992) An introduction to kernel and nearest-neighbor non-parametric regression. *Am Stat* 46(3):175–185
- Arce-Medina E, Paz-Paredes JI (2009) Artificial neural network modeling techniques applied to the hydrodesulfurization process. *Math Comput Model* 49(1–2):207–214
- Bagheri M, Ibrahim ZZ, Mansor SB, Manaf LA, Badarulzaman N, Vaghefi N (2019) Shoreline change analysis and erosion prediction using historical data of Kuala Terengganu, Malaysia. *Environ Earth Sci* 78(15):477
- Barman NK, Chatterjee S, Khan A et al (2014) Trends of shoreline position: an approach to future prediction for Balasore shoreline, Odisha, India. *Open J Mar Sci* 5(01):13
- Bazile R, Boucher MA, Perreault L, Leconte R (2017) Verification of ECMWF system 4 for seasonal hydrological forecasting in a northern climate. *Hydrol Earth Syst Sci* 21(11):5747
- Bosilovich MG, Chen J, Robertson FR, Adler RF (2008) Evaluation of global precipitation in reanalyses. *J Appl Meteorol Climatol* 47(9):2279–2299
- Bosilovich MG, Robertson FR, Takacs L, Molod A, Mocko D (2017) Atmospheric water balance and variability in the MERRA-2 reanalysis. *J Clim* 30(4):1177–1196
- Bouguerra H, Tachi SE, Derdous O, Bouanani A, Khanchoul K (2019) Suspended sediment discharge modeling during flood events using two different artificial neural network algorithms. *Acta Geophys* 67(6):1649–1660
- Bruun P (1962) Sea-level rise as a cause of shore erosion. *J Waterw Harb Div* 88(1):117–132
- Canny JF (1986) A theory of edge detection. *IEEE Trans Pattern Anal Mach Intell* 8:147–163
- Chalabi A, Mohd-Lokman H, Mohd-Suffian I, Karamali K, Karthikeyan V, Masita M (2006) Monitoring shoreline change using ikonos image and aerial photographs: a case study of kuala terengganu area, Malaysia. In: ISPRS Commission VII mid-term symposium “Remote sensing: from pixels to processes”, Enschede, The Netherlands, pp 8–11
- Chudzian P (2011) Radial basis function kernel optimization for pattern classification. In: Burduk R, Kurzyński M, Woźniak M, Żołnierek A (eds) *Computer recognition systems*, vol 4. Springer, Berlin, pp 99–108
- Coltori M (1997) Human impact in the holocene fluvial and coastal evolution of the Marche region, central Italy. *Catena* 30(4):311–335

- Cortes C, Vapnik V (1995) Support-vector networks. *Mach Learn* 20(3):273–297
- Dada OA, Agbaje AO, Adesina RB, Asiwaju-Bello YA (2019) Effect of coastal land use change on coastline dynamics along the Nigerian Transgressive Mahin mud coast. *Ocean Coast Manag* 168:251–264
- De Jong SM, Van der Meer FD (2007) Remote sensing image analysis: including the spatial domain, vol 5. Springer, Berlin
- de Rosnay P, Munoz-Sabater J, Albergel C, Isaksen L, English S, Drusch M, Wigneron JP (2020) SMOS brightness temperature forward modelling and long term monitoring at ECMWF. *Remote Sens Environ* 237(111):424
- Dee DP, Uppala SM, Simmons A, Berrisford P, Poli P, Kobayashi S, Andrae U, Balmaseda M, Balsamo G, Bauer DP et al (2011) The era-interim reanalysis: configuration and performance of the data assimilation system. *Q J R Meteorol Soc* 137(656):553–597
- Dellepiane S, De Laurentiis R, Giordano F (2004) Coastline extraction from sar images and a method for the evaluation of the coastline precision. *Pattern Recogn Lett* 25(13):1461–1470
- Di Silvio G, Nones M (2014) Morphodynamic reaction of a schematic river to sediment input changes: analytical approaches. *Geomorphology* 215:74–82
- Dickens K, Armstrong A (2019) Application of machine learning in satellite derived bathymetry and coastline detection. *SMU Data Sci Rev* 2(1):1–25
- Dolan R, Fenster MS, Holme SJ (1991) Temporal analysis of shoreline recession and accretion. *J Coast Res* 7:723–744
- Dutta D, Mandal A, Afzal MS (2020) Discharge performance of plan view of multi-cycle w-form and circular arc labyrinth weir using machine learning. *Flow Meas Instrum* 73:101740
- ECMWF (2018) European centre for medium-range weather forecasts. <https://www.ecmwf.int/en/research/modelling-and-prediction/marine>
- Elko N, Sallenger A, Guy K, Stockdon H, Morgan K (2002) Barrier island elevations relevant to potential storm impacts: 1. Techniques. US Geological Survey Open File Report, pp 02–287
- Estevés LS, Williams JJ, Dillenburg SR (2006) Seasonal and inter-annual influences on the patterns of shoreline changes in Rio Grande do Sul, southern Brazil. *J Coast Res* 22:1076–1093
- Fadel S, Ghoniemy S, Abdallah M, Sorra HA, Ashour A, Ansary A (2016) Investigating the effect of different kernel functions on the performance of SVM for recognizing Arabic characters. *Int J Adv Comput Sci Appl* 7(1):446–450
- Garg A, Huang H, Kushvaha V, Madhushri P, Kamchoom V, Wani I, Koshy N, Zhu HH (2019) Mechanism of biochar soil pore–gas–water interaction: gas properties of biochar-amended sandy soil at different degrees of compaction using knn modeling. *Acta Geophys* 68:207–217
- Gatys LA, Ecker AS, Bethge M (2015) A neural algorithm of artistic style. [arXiv:150806576](https://arxiv.org/abs/1508.06576)
- Gazi AH, Afzal MS (2020) A new mathematical model to calculate the equilibrium scour depth around a pier. *Acta Geophys* 68(1):181–187
- Gazi AH, Afzal MS, Dey S (2019) Scour around piers under waves: current status of research and its future prospect. *Water* 11(11):2212
- Gelaro R, McCarty W, Molod A, Suarez M, Takacs L, Todling R (2014) The NASA modern era reanalysis for research and applications, Version-2 (MERRA-2). AGU FM 2014:NG32A–01
- Gelaro R, McCarty W, Suárez MJ, Todling R, Molod A, Takacs L, Randles CA, Darmenov A, Bosilovich MG, Reichle R et al (2017) The modern-era retrospective analysis for research and applications, version 2 (merra-2). *J Clim* 30(14):5419–5454
- Govindaraju RS (2000) Artificial neural networks in hydrology. i: preliminary concepts. *J Hydrol Eng* 5(2):115–123. [https://doi.org/10.1061/\(ASCE\)1084-0699\(2000\)5:2\(115\)](https://doi.org/10.1061/(ASCE)1084-0699(2000)5:2(115))
- Govindaraju RS (2000) Artificial neural networks in hydrology. ii: hydrologic applications. *J Hydrol Eng* 5(2):124–137. [https://doi.org/10.1061/\(ASCE\)1084-0699\(2000\)5:2\(124\)](https://doi.org/10.1061/(ASCE)1084-0699(2000)5:2(124))
- Green B (2002) Canny edge detection tutorial. Retrieved 6 Mar 2005
- Gregory K (2004) River channel management. Hodder Education, London
- Guerrero M, Latosinski F, Nones M, Szupiany RN, Re M, Gaeta MG (2015) A sediment fluxes investigation for the 2-d modelling of large river morphodynamics. *Adv Water Resour* 81:186–198
- Gunawardena Y, Ilic S, Pinkerton H, Romanowicz R (2009) Nonlinear transfer function modelling of beach morphology at Duck, North Carolina. *Coast Eng* 56(1):46–58
- Gunn SR et al (1998) Support vector machines for classification and regression. *ISIS Tech Rep* 14(1):5–16
- Halpern BS, McLeod KL, Rosenberg AA, Crowder LB (2008) Managing for cumulative impacts in ecosystem-based management through ocean zoning. *Ocean Coast Manag* 51(3):203–211
- Harley MD, Kinsela MA, Sánchez-García E, Vos K (2019) Shoreline change mapping using crowd-sourced smartphone images. *Coast Eng* 150:175–189
- Hashemi M, Ghadampour Z, Neill S (2010) Using an artificial neural network to model seasonal changes in beach profiles. *Ocean Eng* 37(14–15):1345–1356
- Houser C, Hapke C, Hamilton S (2008) Controls on coastal dune morphology, shoreline erosion and barrier island response to extreme storms. *Geomorphology* 100(3–4):223–240
- Howarth PJ, Wickware GM (1981) Procedures for change detection using landsat digital data. *Int J Remote Sens* 2(3):277–291
- Hsu HH, Hoskins BJ (1989) Tidal fluctuations as seen in ECMWF data. *Q J R Meteorol Soc* 115(486):247–264
- Hsu CW, Chang CC, Lin CJ et al (2003) A practical guide to support vector classification. Department of Computer Science National Taiwan University
- Hu LY, Huang MW, Ke SW, Tsai CF (2016) The distance function effect on k-nearest neighbor classification for medical datasets. *SpringerPlus* 5(1):1304
- Jan J, Hung SL, Chi S, Chern J (2002) Neural network forecast model in deep excavation. *J Comput Civ Eng* 16(1):59–65
- Jangir B, Satyanarayana A, Swati S, Jayaram C, Chowdary V, Dadhwal V (2016) Delineation of spatio-temporal changes of shoreline and geomorphological features of Odisha coast of India using remote sensing and gis techniques. *Nat Hazards* 82(3):1437–1455
- Kennedy AD, Dong X, Xi B, Xie S, Zhang Y, Chen J (2011) A comparison of MERRA and NARR reanalyses with the DOE ARM SGP data. *J Clim* 24(17):4541–4557
- Kesikoğlu MH, Çiçekli SY, Kaynak T (2020) The identification of coastline changes from landsat 8 satellite data using artificial neural networks and K-nearest neighbor. *Turk J Eng* 4(1):47–56
- Khaledian M, Isazadeh M, Biazar S, Pham Q (2020) Simulating Caspian sea surface water level by artificial neural network and support vector machine models. *Acta Geophys* 68:553–563
- Kim IH, Lee HS, Song DS (2013) Time series analysis of shoreline changes in Gonghyunjin and Songjiho Beaches, South Korea using aerial photographs and remotely sensed imagery. *J Coast Res* 65:1415–1420
- Kumar TS, Mahendra R, Nayak S, Radhakrishnan K, Sahu K (2010) Coastal vulnerability assessment for Orissa State, east coast of India. *J Coast Res* 26:523–534
- Larson M, Capobianco M, Hanson H (2000) Relationship between beach profiles and waves at Duck, North Carolina, determined by canonical correlation analysis. *Mar Geol* 163(1–4):275–288



- Lee YK, Eom J, Do JD, Kim BJ, Ryu JH (2019) Shoreline movement monitoring and geomorphologic changes of beaches using Lidar and UAVs Images on the Coast of the East Sea, Korea. *J Coast Res* 90(sp1):409–414
- Li R, Liu JK, Felus Y (2001) Spatial modeling and analysis for shoreline change detection and coastal erosion monitoring. *Mar Geod* 24(1):1–12
- Markose VJ, Rajan B, Kankara R, Selvan SC, Dhanalakshmi S (2016) Quantitative analysis of temporal variations on shoreline change pattern along Ganjam district, Odisha, East Coast of India. *Environ Earth Sci* 75(10):929
- MERRA-2 (2017) Modern era retrospective-analysis for research and applications. <https://gmao.gsfc.nasa.gov/reanalysis/MERRA-2/>
- Mishra M, Chand P, Pattnaik N, Kattel DB, Panda G, Mohanti M, Baruah UD, Chandniha SK, Achary S, Mohanty T (2019) Response of long-to short-term changes of the Puri coastline of Odisha (India) to natural and anthropogenic factors: a remote sensing and statistical assessment. *Environ Earth Sci* 78(11):338
- Monalisha M, Panda G (2018) Coastal erosion and shoreline change in Ganjam coast along East Coast of India. *J Earth Sci Clim Change* 9:467
- Montaño J, Coco G, Antolínez JA, Beuzen T, Bryan KR, Cagigal L, Castelle B, Davidson MA, Goldstein EB, Ibaceta R et al (2020) Blind testing of shoreline evolution models. *Sci Rep* 10(1):1–10
- Morton R (1996) Geoinicators of coastal wetlands and shorelines. Geoinicators: assessment rapid environmental changes in earth systems. AA Balkema, Rotterdam, pp 207–230
- Mukhopadhyay A, Mukherjee S, Mukherjee S, Ghosh S, Hazra S, Mitra D (2012) Automatic shoreline detection and future prediction: a case study on Puri Coast, Bay of Bengal, India. *Eur J Remote Sens* 45(1):201–213
- Murthy VS, Gupta S, Mohanta D (2009) Distribution system insulator monitoring using video surveillance and support vector machines for complex background images. *Int J Power Energy Convers* 1(1):49–72
- Nandi S, Ghosh M, Kundu A, Dutta D, Baksi M (2016) Shoreline shifting and its prediction using remote sensing and gis techniques: a case study of Sagar Island, West Bengal (India). *J Coast Conserv* 20(1):61–80
- Nowakowski A (2015) Remote sensing data binary classification using boosting with simple classifiers. *Acta Geophys* 63(5):1447–1462
- Peponi A, Morgado P, Trindade J (2019) Combining artificial neural networks and gis fundamentals for coastal erosion prediction modeling. *Sustainability* 11(4):975
- Pescaroli G, Nones M, Galbusera L, Alexander D (2018) Understanding and mitigating cascading crises in the global interconnected system. *Int J Disaster Risk Reduction* 30:159–163
- Piasecki A, Jurasz J, Adamowski JF (2018) Forecasting surface water-level fluctuations of a small glacial lake in Poland using a wavelet-based artificial intelligence method. *Acta Geophys* 66(5):1093–1107
- Pierini JO, Lovallo M, Telesca L, Gómez EA (2013) Investigating prediction performance of an artificial neural network and a numerical model of the tidal signal at Puerto Belgrano, Bahía Blanca Estuary (Argentina). *Acta Geophys* 61(6):1522–1537
- Puskarczyk E (2019) Artificial neural networks as a tool for pattern recognition and electrofacies analysis in Polish palaeozoic shale gas formations. *Acta Geophys* 67(6):1991–2003
- Rajawat A, Chauhan H, Ratheesh R, Rode S, Bhandari R, Mahapatra M, Kumar M, Yadav R, Abraham S, Singh S et al (2015) Assessment of coastal erosion along the Indian Coast on 1: 25,000 scale using satellite data of 1989–1991 and 2004–2006 time frames. *Curr Sci* 109:347–353
- Ramesh R, Purvaja R, Senthil Vel A (2011) National assessment of shoreline change: Odisha coast. NCSCM/ MoEF Report 2011-01, 57 p., available at <http://www.ncscm.org/reports.php>
- Ramesh R, R P, Vel S (2017) A shoreline change assessment for Odisha Coast; National Centre for Sustainable Coastal Management (NCSCM). Govt. of Odisha Report. National Centre for Sustainable Coastal Management (NCSCM). Accessed on 11 Nov 2017
- Reichle RH, Koster RD, De Lannoy GJ, Forman BA, Liu Q, Mahanama SP, Touré A (2011) Assessment and enhancement of merra land surface hydrology estimates. *J Clim* 24(24):6322–6338
- Rienecker MM, Suarez MJ, Gelaro R, Todling R, Bacmeister J, Liu E, Bosilovich MG, Schubert SD, Takacs L, Kim GK et al (2011) Merra: Nasa's modern-era retrospective analysis for research and applications. *J Clim* 24(14):3624–3648
- Ronco P, Fasolato G, Nones M, Di Silvio G (2010) Morphological effects of damming on lower Zambezi river. *Geomorphology* 115(1–2):43–55
- Ryan T, Sementilli P, Yuen P, Hunt B (1991) Extraction of shoreline features by neural nets and image processing. *Photogramm Eng Remote Sens* 57(7):947–955
- Saluja S, Singh AK, Agrawal S (2013) A study of edge-detection methods. *Int J Adv Res Comput Commun Eng* 2(1):994–999
- Satapathy SC, Udgata SK, Biswal BN (2012) Proceedings of the international conference on frontiers of intelligent computing: theory and applications (FICTA), vol 199. Springer, Berlin
- Schalkoff RJ (1997) Artificial neural networks, vol 1. McGraw-Hill, New York
- Shen S, Ostrenga D, Vollmer B, Li A, Meyer D (2019) MERRA-2 data and analytic services at NASA GES DISC for climate extremes study. In: 16th AOGS-Annual meeting of asia oceania geosciences society, July 28, 2019–August 02, 2019, Singapore
- Shen S, Ostrenga DM, Bosilovich MG, Li AW, Meyer DJ (2020) Near 40 years MERRA-2 data at NASA GES DISC-opportunity and challenge to support extremes study. In: 100th AMS Annual Meeting, January 12, 2020–January 16, 2020, Boston, United States
- Shrivakshan G, Chandrasekar C (2012) A comparison of various edge detection techniques used in image processing. *Int J Comput Sci Issues: IJCSI* 9(5):269
- Simeoni U, Corbau C (2009) A review of the delta po evolution (Italy) related to climatic changes and human impacts. *Geomorphology* 107(1–2):64–71
- Small C, Nicholls RJ (2003) A global analysis of human settlement in coastal zones. *J Coast Res* 19:584–599
- Sobel I, Feldman G (1968) A  $3 \times 3$  isotropic gradient operator for image processing. A talk at the Stanford artificial project, pp 271–272
- Stockdon HF, Doran KS, Sallenger AH Jr (2009) Extraction of lidar-based dune-crest elevations for use in examining the vulnerability of beaches to inundation during hurricanes. *J Coast Res* 53:59–65
- Suanez S, Cariolet JM, Cancouët R, Arduin F, Delacourt C (2012) Dune recovery after storm erosion on a high-energy beach: Vougot Beach, Brittany (France). *Geomorphology* 139:16–33
- The Indian Tide Tables-Part 1,1995: Indian and Selected Foreign Ports (1994) Surveyor general of India, printed by survey of India, Dehradun
- Tsekouras GE, Trygonis V, Maniatopoulos A, Rigos A, Chatzispavlis A, Tsimikas J, Mitianoudis N, Velegrakis AF (2018) A hermite neural network incorporating artificial bee colony optimization to model shoreline realignment at a reef-fronted beach. *Neurocomputing* 280:32–45
- USGS (2017) United states geological survey. <https://earthexplorer.usgs.gov>
- Valiela I (2004) Global coastal change. Blackwell, Oxford

- Valipour M, Tian D (2018) Comparing soil moisture dynamics in climate reanalyses, land surface models, and remote sensing retrievals over the continental united states. In: AGU Fall Meeting Abstracts
- Valipour M, Banihabib M, Behbahani S (2012) Monthly inflow forecasting using autoregressive artificial neural network. *J Appl Sci* 12(20):2139–2147
- Valipour M, Banihabib ME, Behbahani SMR (2013) Comparison of the arma, arima, and the autoregressive artificial neural network models in forecasting the monthly inflow of dez dam reservoir. *J Hydrol* 476:433–441
- Vapnik V (1963) Pattern recognition using generalized portrait method. *Autom Remote Control* 24:774–780
- Vapnik VN, Chervone AY (1965) On a class of pattern-recognition learning algorithms. *Autom Remote Control* 25(6):838
- Varrani A, Nones M, Gupana R (2019) Long-term modelling of fluvial systems at the watershed scale: examples from three case studies. *J Hydrol* 574:1042–1052
- Vijayarani S, Vinupriya M (2013) Performance analysis of Canny and Sobel edge detection algorithms in image mining. *Int J Innov Res Comput Commun Eng* 1(8):1760–1767
- Vincent OR, Folorunso O et al (2009) A descriptive algorithm for sobel image edge detection. In: Proceedings of informing science & IT education conference (InSITE), vol 40. Informing Science Institute California, pp 97–107
- Wang J, Li B, Gao Z, Wang J (2019) Comparison of ECMWF significant wave height forecasts in the China sea with buoy data. *Weather Forecast* 34(6):1693–1704
- White K, El Asmar HM (1999) Monitoring changing position of coastlines using Thematic Mapper imagery, an example from the Nile Delta. *Geomorphology* 29(1–2):93–105
- Zhang X, Wang Z (2010) Coastline extraction from remote sensing image based on improved minimum filter. In: 2010 second IITA international conference on geoscience and remote sensing, vol 2. IEEE, pp 44–47



# Dynamics of thin disk settling in two-layered fluid with density transition

Magdalena M. Mrokowska<sup>1</sup>

Received: 27 March 2020 / Accepted: 13 June 2020 / Published online: 25 June 2020  
© The Author(s) 2020

## Abstract

Settling of solid particles in a stratified ambient fluid is a process widely encountered in geophysical flows. A set of experiments demonstrating the settling behaviour (the pattern of trajectory, variation of particle orientation, and settling velocity with depth) of thin disks descending through a nonlinear density transition was performed. The results showed complex hydrodynamic interactions between a particle and a liquid causing settling orientation instabilities and unsteady particle descent in low to moderate Reynolds number regime. Five phases of settling were observed: two phases with stable horizontal, one with stable vertical disk position, and two reorientation phases; moreover, two local minima of settling velocity were identified. It was demonstrated that thresholds for local minima and the first reorientation depend on the settling dynamics in an upper layer, stratification conditions, and disk geometry. The comparison of settling behaviour of thin disks varying in diameter revealed that settling dynamics is sensitive to particle geometry mainly in the upper part of density transition with a non-obvious result that the first minimum velocity is smaller for a disk with a larger diameter than for a disk with a smaller diameter. The analysis of settling trajectory showed that two reorientations are accompanied with a horizontal drift, which may be important in the context of interactions between particles settling in a group.

**Keywords** Stratification · Particle settling · Disk · Density transition

## Introduction

Density stratification occurs in various fluid components of natural environment (ocean, atmosphere, and the Earth's interior) and affects to a large extent the vertical transport of particles. Modification of settling or rising behaviour of rigid particles, drops, and bubbles due to the presence of sharp or continuous stratification may considerably influence geophysical processes.

Density gradients form in aquatic systems (ocean, seas, and lakes) as a consequence of temperature and/or salinity variation with depth. The settling dynamics of particles in a density-stratified ambient is much different from that in homogeneous conditions. Sharp density gradients known as pycnoclines (haloclines or thermoclines, with salinity or temperature acting as a stratifying agent, respectively) form in favourable conditions (Capet et al. 2016; Noufal

et al. 2017) induced the deceleration and prolonged residence times of particles in the stratified region (Peperzak et al. 2003). Field observations have provided evidence that organic particles such as marine, lake snow, and faecal pellets may stagnate at pycnoclines for a few days, forming the so-called thin layers (Diercks et al. 2019; Macintyre et al. 1995; Prairie et al. 2015), which modify particulate organic carbon flux in the ocean (Arnosti 2011; Lutz et al. 2002; Prairie et al. 2017).

Atmosphere stratification affects the transport dynamics of various particles including dust, aerosol, pollens, volcanic ashes, and pollutants. Temperature inversion layer that may form in the troposphere prevents airborne particles from settling, which has been observed for dust (Zhai et al. 2019). Moreover, the stratification of atmosphere affects the fate of eruption columns and the transport of volcanic particles including settling of ashes (Woods 1995). In the context of Earth's interior, it has been demonstrated that the presence of sinking crystals may considerably accelerate mixing between rhyolitic and basaltic magmas, suggesting that particle settling should be considered in the magma mixing process (Renggli et al. 2016).

✉ Magdalena M. Mrokowska  
m.mrokowska@igf.edu.pl

<sup>1</sup> Institute of Geophysics, Polish Academy of Sciences, Ks. Janusza 64, 01-452 Warsaw, Poland

The effects of fluid stratification on the movement of an object are also important in oceanographic techniques allowing to measure physical, e.g. temperature and salinity, as well as biological parameters across the ocean depth. The design and operation of oceanographic float take advantage of natural density stratification, since the device actively changes its depth by manipulating its density with reference to ambient conditions to achieve the equilibrium depth. Although the dimensions of floats considerably exceed the above-mentioned natural particles, the same physical processes govern the settling and rising of objects in stratified fluid (D'Asaro 2018).

Research on the dynamics of particle settling in stratified conditions is still at the level of fundamental mechanics of particle motion, since density gradient effects considerably increase the complexity of the problem compared to the homogeneous conditions. While settling in a homogeneous fluid is conveniently characterised by Reynolds number,  $Re = U a/\nu$ , where  $U$  is the settling velocity [ $\text{m s}^{-1}$ ],  $a$ —particle characteristic length [m],  $\nu$ —kinematic viscosity [ $\text{m}^2 \text{s}^{-1}$ ], when stratified fluid is considered, stratification effects are usually accounted by stratification strength expressed as the Brunt–Vaisala buoyancy frequency  $N$ , and Froude number,  $Fr = U/N a$ , which is the ratio of inertial to buoyancy forces (Yick et al. 2009). Another parameter is the ratio between momentum diffusivity and mass diffusivity of stratifying agent, i.e. Schmidt number,  $Sc = \nu/\kappa$  where  $\kappa$  is the diffusivity of stratifying agent [ $\text{m}^2 \text{s}^{-1}$ ] (or Prandtl number for temperature).

Since detailed research on individual particle dynamics is challenging in natural conditions, small-scale laboratory experiments and numerical studies have become the major means to extend our knowledge on the fundamental aspects of sedimentation process (Prairie and White 2017). Earlier studies considered linear ambient stratification (Doostmohammadi et al. 2014; Mercier et al. 2020; Yick et al. 2009), a two-layered configuration with a sharp density transition (Abaid et al. 2004; Camassa et al. 2010; Srdic-Mitrovic et al. 1999; Verso et al. 2019), and a two-layered configuration with a continuous nonlinear transition (Mrokowska 2018). Previous experimental studies have referred to the sedimentation in stratified conditions in a general sense, however, with a strong focus on natural waters where stratification is generated by the vertical variation of salinity and temperature.

Density of particles present in aquatic systems (detrital material, mineral particles, plankton, microplastics, and marine snow) is close to that of water, which in combination with small dimensions of particles makes them settle in a low and moderate Reynolds number regime. Viscous forces dominate for  $Re \ll 1$ , while inertial forces affect the settling dynamics for  $Re$  higher than unity. A group of studies motivated by settling processes in marine systems focused

on settling in the viscous regime (Camassa et al. 2010; Yick et al. 2009), while others considered low to moderate  $Re$  number inertia-controlled regimes (Kindler et al. 2010; Mrokowska 2018; Prairie et al. 2015).

It has been well acknowledged in research performed so far that pronounced deceleration of particles at pycnoclines observed in nature is due to stratification-induced drag, which appears in the presence of density gradient beside the drag characteristic for homogeneous conditions (Magnaudet and Mercier 2020; Srdic-Mitrovic et al. 1999). The origin of stratification-induced drag is attributed to a caudal fluid entrained in the wake of particle from the above layers of lighter fluid (Srdic-Mitrovic et al. 1999) and to the compression and distortion of isopycnals (Doostmohammadi et al. 2012; Yick et al. 2009). These two basic mechanisms have been demonstrated for spheres using both laboratory experiments and numerical simulations. Inertial waves generated by descending particle may be the source of additional drag (Srdic-Mitrovic et al. 1999; Yick et al. 2009); however, this applies only for moderate and high  $Re$  number regimes in which case the particle has large inertia (Okino et al. 2017; Scase and Dalziel 2004).

Although the knowledge on the mechanisms of particle deceleration in the presence of stratification is growing, it is still not sufficient to propose robust methods that could be applied in sedimentation studies by a wide community of earth and environmental sciences researchers. Existing sedimentation and biogeochemical models oversimplify ambient conditions due to insufficient existing knowledge on how to tackle complex physical properties of natural waters in combination with settling processes, which may cause misestimating of sedimentation fluxes (Lutz et al. 2002). Effects of stratification are either discarded or limited to the effect of density change with the assumption that formulas for homogeneous conditions hold in stratified ambient, while settling velocity is much lower in a stratified background configuration than predicted by the standard Stokes law formula (Dey et al. 2019) for the corresponding homogeneous conditions due to stratification-induced drag (Camassa et al. 2009), which has serious implications on the estimation of particulate flux. Despite some successful attempts to parametrize the sedimentation through a pycnocline (Prairie and White 2017), research on sedimentation in stratified conditions has not yet provided robust methods to be applied easily in large-scale or local-scale models.

Another troublesome factor affecting settling process is the shape of a particle. The vast majority of particles present in natural fluid systems are non-spherical; some examples are irregular volcanic ashes (Saxby et al. 2018), aeolian sediments (Raffaele et al. 2020), microplastics (Cole et al. 2011), faecal pellets, microorganisms, and marine snow (Turner 2015). Basic research within fluid mechanics and sedimentology provides some theoretical bases for the estimation

of settling velocities (Dietrich 1982; Loth 2008); however, existing theoretical background for the assessment of settling dynamics of solid particles of various shapes and densities is not sufficient to effectively study geophysical and environmental problems such as settling speeds and fluxes of marine snow (Laurenceau-Cornec et al. 2019), microplastics in the ocean (Khatmullina and Isachenko 2017), and particulate organic matter flux (Lutz et al. 2002). The major problem is the fact that settling dynamics of variously shaped particles cannot be described by a universal law, as in the case of a sphere, and attempts to apply formulae dedicated for a sphere results in unreliable settling velocity estimations (Saxby et al. 2018). Consequently, an excessive number of semi-empirical relations to calculate drag coefficient and settling velocity have been derived. However, these relations have serious limitations on their applicability due to shape-dependent and orientation-dependent drag (Bagheri and Bonadonna 2016; Loth 2008). These practical problems have been already faced in research on settling fluxes of microplastics (Khatmullina and Isachenko 2017; Waldschlager and Schuttrumpf 2019) and biogenic particles (Maggi 2013).

Since our state of knowledge is not sufficient to properly account for stratification effects in sedimentation models, fundamental studies on settling dynamics of non-spherical particles in stratified ambient are necessary to formulate in the future effective methods that will be applied in models. A few basic studies performed so far demonstrated that shape effects have pronounced impact on particle settling dynamics in stratified systems (Doostmohammadi and Ardekani 2014; Mercier et al. 2020; Mrokowska 2018). Spheroid particles reorient at a density interface as the effect of buoyancy-induced torque appearing due to pressure difference at the edges of particle which overcomes inertial torque, which leads to complex pattern of settling behaviour (Ardekani et al. 2017; Doostmohammadi and Ardekani 2014; Mrokowska 2018). Disk rotation has been claimed to be additionally amplified by the torque exerted by a jet moving from the centre of the disk to its edge (Mercier et al. 2020), an effect visualized also in Mrokowska (2018).

Disk-like particles constitute a large group of particles that are of significance in environmental processes, e.g. diatoms, minerals, and ashes, and are convenient representation of oblate spheroids. Therefore, their dynamics are worth studying in the context of sedimentation in natural fluid systems. The pattern of individual disk settling has been described for two ambient conditions so far: a linear (Mercier et al. 2020) and nonlinear stratification (Mrokowska 2018). In the first study, short-length cylinders settling in a linearly stratified liquid have been considered (Mercier et al. 2020), and in the second one thin disks settling through a density transition have been studied (Mrokowska 2018). In the linear stratification set-up, the whole water column has been stratified with the maximum density exceeding the

density of a particle, while in the other case there has been a two-layered set-up with two homogeneous layers of liquid: a less dense upper and a denser lower, both of density smaller than the density of the particle, with a nonlinear density transition between the layers. For both configurations, a disk reorients from stable horizontal to stable vertical position after some distance fallen broadside on in a stratified layer, which has been described as phase 1 and 2 in Mercier et al. (2020) and phases 1, 2, 3 in Mrokowska (2018). However, the dynamics of a disk settling in a nonlinear stratification are more complex, which manifests in two local velocity minima, while no minima have been observed for linear stratification. It should be noted that the first velocity minimum is an effect of particle deceleration during reorientation and the second one is achieved when a disk is in vertical position, but occurs only when stratification is sufficiently strong (Mrokowska 2018). Two different mechanisms have been described for the reorientation from vertical to horizontal position. In the case of a linear stratification, reorientation occurs when the particle achieves the level of neutral buoyancy, which is described as phase 3 in (Mercier et al. 2020), while in the nonlinear stratification the second reorientation is an effect of fading influence of stratification when the particle enters the lower homogeneous layer where the role of inertia is restored. This is described as phase 4, and further settling in a stable horizontal position in the lower layer is denoted as phase 5 (Mrokowska 2018).

Another study (Lam et al. 2019) has investigated the settling of heavy disks in a stably stratified liquid column for  $Re$  of order 1000. In this study, the effect of stratification on secondary motion of disk, specifically a fluttering mode, has been examined. The study showed that, compared to a homogeneous fluid, stratification modifies settling dynamics enhancing particle radial dispersion, decreasing settling velocity, inclination angle, and uttering amplitude.

The settling process through stratified fluid has not been yet well recognised at the fundamental level, and deeper insight into settling dynamics of non-spherical particles in a stratified fluid is necessary to enable development of methods that will improve the performance of sedimentation estimations and particle flux modelling in future. This study aims at improving our understanding of non-spherical particles settling in stratified conditions. I show a series of settling experiments in a two-layered water column with nonlinear density transition extending the previous research (Mrokowska 2018) to focus on new aspects of settling dynamics. To gain more detailed characteristics of disk settling behaviour in a transition layer, 3D trajectory of particle settling was retrieved. The settling dynamics of thin disks were evaluated for a range of settling conditions in an upper layer and various stratification strengths in a density transition region. Two types of thin disks varying in diameter were used in this study. The objectives of the

study was to (1) quantify the effect of settling conditions in an upper layer and stratification characteristics on the settling behaviour, i.e. settling velocity, reorientation pattern, and particle trajectory of disks in a low and moderate Re number regime, and to (2) assess the effect of disk diameter on settling dynamics.

## Methods

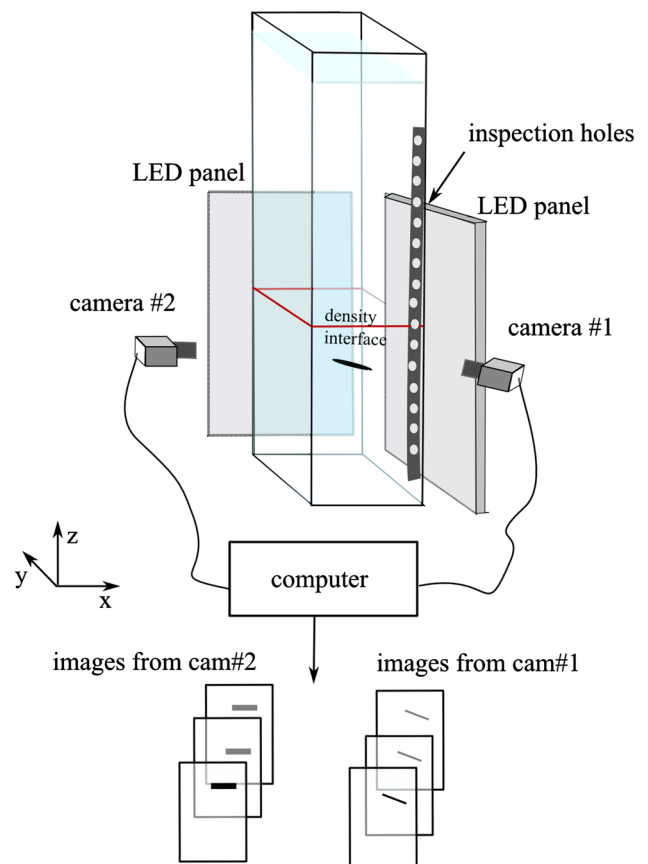
### Characterization of particles and experimental set-up

Experiments were carried out in the Laboratory of Hydrodynamic Micromodels, Institute of Geophysics, Polish Academy of Sciences, Warsaw, Poland. Disks were manufactured from acrylonitrile butadiene styrene (ABS) foil with a density of  $1050 \pm 5 \text{ kg m}^{-3}$  and thickness  $h = 50 \mu\text{m}$ . Specially designed puncher was used to produce two sets of disks differing in diameter,  $d = 2 \text{ mm}$  and  $3 \text{ mm}$  with  $0.1 \text{ mm}$  accuracy, denoted as  $d2$  and  $d3$  disks, respectively. Disks were considered thin with a geometrical aspect ratio  $\chi = d/h$  equal to 40 and 60. Several disks of each type were prepared, and diameters were verified by image analysis using one of cameras applied in this study to take high-resolution photographs.

Disk settling was investigated in a series of experiments carried out in a 0.50-m-high transparent tank with a square base ( $0.10 \times 0.10 \text{ m}$ ). A two-layered water column was formed before an experiment by filling the tank with denser salt water solution (density  $\rho_{ll}$ ) up to about 0.19 m from the bottom and with less dense salt water solution (density  $\rho_{ul}$ ) from 0.19 m up to the height of about 0.48 m. Thereby, two layers referred to as an upper and lower layer were formed (Fig. 1). To fill the tank, a method applied in the previous study has been used where details are described (Mrokowska 2018).

To estimate vertical variation of salinity, a procedure described in the mentioned paper was applied. Inspection holes spaced every 5 mm (Fig. 1) were used to sample aqueous salt solution, and the salinity of each sample was measured using Kruss refractometer, model DR301-95. Given the temperature measurements, salinity was recalculated to density using literature tables (Kestin et al. 1981). The vertical variation of density was fitted to the hyperbolic tangent function (Prairie et al. 2015), as in previous study (Mrokowska 2018). The experiments were carried out at room temperature varying in a range  $22.1\text{--}23.4 \text{ }^\circ\text{C}$ .

A particle was released beneath a water surface in the centre of  $x$ - $y$  plane, and particle descent in density transition region was recorded using two identical cameras (Basler acA2500-60um equipped with Schneider-Kreuznach macro lenses Componon 2.8/28-001) positioned orthogonal to each



**Fig. 1** Experimental set-up showing a two-camera configuration for measuring 3D disk trajectory

other. Camera#1 filmed  $x$ - $z$  plane and camera#2  $y$ - $z$  plane (Fig. 1), which enabled the reconstruction of particle trajectory. Two LED panels with DC power supply were used to generate backlight necessary to visualise the settling particle and to record the projection of particle shadow. The field of view (FOV) of a camera covered 77-mm-high and 62-mm-wide area with one pixel corresponding to  $31 \mu\text{m}$ . Upper edges of both cameras FOV were precisely positioned to the same level slightly above the density interface. Particle trajectories were recorded at 60 fps, and the capturing of image pairs was synchronized. Each settling test was repeated a few times.

The settling of a few disks of the same type was investigated in each experiment to check the repeatability of an experiment. The number of repetitions varied from 2 to 5 and was constrained by (1) the quality of visualisation in two orthogonal planes (tests were discarded from analysis when a particle was settling not in the centreline of the tank) and (2) timescale of mixing in transition layer due to diffusion and mixing induced by settling particles.

Experimental set-ups were checked for the optical distortions, due to possible change of refractive index in

density-stratified fluid, using a ruler (Abaid et al. 2004) and examining the dimensions of settling particles in experimental images; no deformations have been observed.

Cameras location enabled recording particle settling only in the transition layer; hence, another set of experiments in homogeneous conditions corresponding to settling conditions in an upper layer were performed in the same tank to assess terminal settling velocity in the upper layer. In these experiments, only one camera was used to record settling particle in FOV positioned about 0.3 m below a free water surface to ensure terminal settling conditions.

### Processing of images

Pairs of orthogonal views of settling disk (images captured by camera#1 and camera#2) were obtained from each experimental test on particle settling within density transition. Image analysis was performed using procedures available in ImageJ and using ad hoc scripts in MATLAB® following methods applied in the previous studies (Mrokowska 2018; Mrokowska and Krztoń-Maziopa 2019) to identify particles and assess their position. Thresholding method was applied to assess the contours of particle projection in each image. Coordinates of particle projection geometrical centre representing the centre of mass were evaluated from image pairs, which were next used to retrieve the particle trajectory. Small vertical shifts (not exceeding 1.1 mm in average) between images in each pair were corrected based on the location of the first minimum velocity (defined further in the test), which is unequivocal in both images. Time-resolved position data were smoothed by a moving-average cubic polynomial using Savitzky–Golay filter. Settling velocity was evaluated as a central-point difference quotient using particle position data and time step between consecutive images (Mrokowska 2018; Mrokowska and Krztoń-Maziopa 2019).

### Experimental conditions

Three sets of experiments in a two-layered configuration were carried out with density of upper layer,  $\rho_{ul}$ , ranging between 1003 and 1016 kg m<sup>-3</sup>, while the density of lower layer was constant ( $\rho_{ll} = 1036$  kg m<sup>-3</sup>). Viscosity of fluid was assumed constant for the purposes of data analysis, since the difference in viscosity between an upper and lower layer was insignificant (smaller than 5%). Names of experiments ES%*d*x reflect salinity in an upper layer and a disk diameter, where *S* stands for salinity [%] and *x* is a diameter [mm]. Details on the physical properties of liquids are reported in Table 1.

Density of ambient fluid varies nonlinearly with depth. Figure 2 shows vertical distribution of ambient fluid density,  $\rho_f$ , both measured and fitted to the hyperbolic tangent function (Eq. 1) with  $R^2 > 0.99$ :

$$\rho_f(z) = \left( \frac{\rho_{ll} - \rho_{ul}}{2} \right) \left( 1 + \tanh \left( \frac{z - z_0}{p} \right) \right) + \rho_{ul} \quad (1)$$

where  $\rho_f$  is the density of fluid [kg m<sup>-3</sup>],  $\rho_{ul}$  and  $\rho_{ll}$ —density of homogeneous upper and lower layer, respectively [kg m<sup>-3</sup>],  $z$ —vertical coordinate [m],  $z_0$ ,  $p$ —fitting parameters [m]. Density jump,  $b = (\rho_{ll} - \rho_{ul})/\rho_{ul}$ , indicating density difference between an upper and lower layer varies between 0.020 and 0.033. Brunt–Vaisala buoyancy frequency,  $N$ , varies within density transition (increases with rising density gradient (Fig. 2)) and was evaluated from the formula:

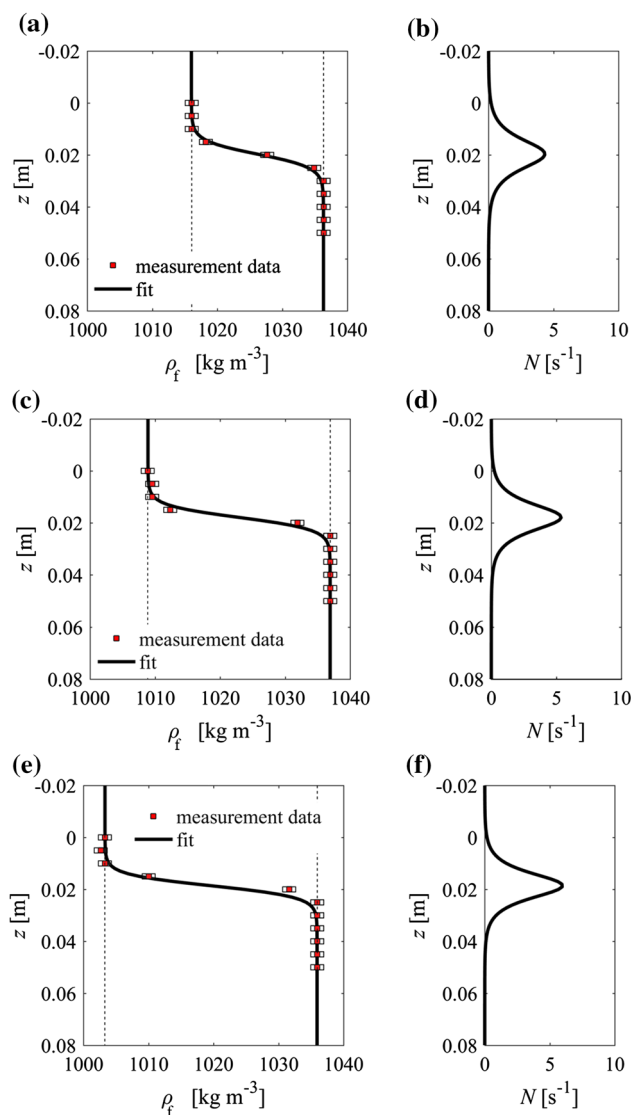
$$N(z) = \sqrt{\frac{g}{\rho_f(z)} \frac{\partial \rho_f(z)}{\partial z}} \quad (2)$$

where  $g = 9.81$ —acceleration due to gravity [m s<sup>-2</sup>],  $\frac{\partial \rho_f(z)}{\partial z}$ —background density gradient. Since there is no specific value of  $N$  as in the linear stratification, the maximum buoyancy frequency  $N_{max}$  is defined to be a general parameter describing stratification strength for the purposes of this

**Table 1** Experimental conditions for two-layered density configuration

Exp	$\rho_{ll}$ [kg m <sup>-3</sup> ]	$\rho_{ul}$ [kg m <sup>-3</sup> ]	$b$ [-]	$\nu \times 10^{-6}$ [m <sup>2</sup> s <sup>-1</sup> ]	$T$ [°C]	$N_{max}$ [s <sup>-1</sup> ]	$L_t/d$ [-]	Total no. of repetitions
E2.6% <i>d</i> 2	1036	1016	0.020	1.0	22.1	4.328	20	4
E2.6% <i>d</i> 3							13	4
E1.6% <i>d</i> 2	1037	1009	0.028	1.0	22.5	5.366	19	3
E1.6% <i>d</i> 3							12	2
E0.7% <i>d</i> 2	1036	1003	0.033	1.0	23.4	5.980	19	3
E0.7% <i>d</i> 3							12	5

$b$ —density jump given as  $(\rho_{ll} - \rho_{ul})/\rho_{ul}$ ,  $N$ —Brunt–Vaisala buoyancy frequency evaluated from Eq. (2),  $N_{max}$ —maximum Brunt–Vaisala buoyancy frequency,  $L_t$ —transition thickness defined as a region where  $N > 0.2$  s<sup>-1</sup>,  $\nu$ —reference kinematic viscosity,  $T$ —temperature of liquid



**Fig. 2** Density profile (measured and fitted using Eq. (1)) for two-layered experiments (a) E2.6%, (c) E1.6%, (e) E0.7% and corresponding variation of Brunt–Vaisala buoyancy frequency,  $N$ , evaluated using Eq. (2),  $z=0$  corresponds with the upper boundary of transition layer (b) E2.6%, (d) E1.6%, (f) E0.7%

study. Buoyancy frequency increases with density jump with  $N_{\max}$  varying from 4.328 to 5.980 (Table 1).

Transition thickness,  $L_t$ , does not vary significantly between the three experiments; however, it decreases slightly with the increasing density jump. Transition thickness relative to disk diameter,  $L_t/d$ , varies from 12 to 20 (Table 1) and is considered continuous, since it exceeds the dimensions of particles in contrast to sharp interface configurations where the thickness of density transition was comparable with particle dimensions (Blanchette and Shapiro 2012; Camassa et al. 2009).

Experiments in homogeneous conditions referring to terminal settling in an upper layer accompanied each two-layered experiment. Two homogeneous experiments, each comprising several repetitions, were performed for each configuration to gather meaningful number of repetitions, and the experimental conditions are reported in Table 2. These data were used to analyse the impact of settling parameters in an upper layer on particle settling behaviour within the transition layer. Reynolds number in an upper layer (or entrance Reynolds number) is defined as  $Re_{ul} = d U_{ul}/\nu$ , where  $U_{ul}$  denotes the terminal settling velocity (Table 2).

This study has been a continuation of previous experimental research (Mrokowska 2018) extending the range of conditions and focusing on the quantification of the effects of density transition on settling behaviour of disks. In the present study, stratification strength was controlled by the density of upper layer varying between three experiments, while the density of lower layer was kept constant. Conversely, the density of upper layer was kept constant in the previous study. Here, the thickness of transition layer,  $L_t$ , was almost constant (0.037–0.040 m) in all experiments. Hence, density jump between upper and lower layer may be considered as the main source of stratification variability within transition layer with negligible impact of transition thickness. This is different configuration than studied in the previous research, where the stratification strength was controlled by the transition thickness. Disks used in the present study were of larger diameter than in the previous one where diameters smaller than 2 mm were investigated.

**Table 2** Experimental conditions for homogeneous fluid configuration referring to an upper layer conditions

Exp.	$\rho_{ul}$ [kg m <sup>-3</sup> ]	$T$ [°C]	$\nu \times 10^{-6}$ [m <sup>2</sup> s <sup>-1</sup> ]	$U_{ul} \pm SD$ [m s <sup>-1</sup> ]	$Re_{ul}$ [-]	$Ar_{ul}$ [-]	Total no. of repetitions
E2.6% $d2$	1016	22.5–23.0	1.0	0.0025 ± 0.0003	5.0	2.9	7
E2.6% $d3$				0.0030 ± 0.0003	9.0	4.4	8
E1.6% $d2$	1009	22.5–23.0	1.0	0.0028 ± 0.0004	5.6	3.4	8
E1.6% $d3$				0.0035 ± 0.0004	10.5	5.1	7
E0.7% $d2$	1003	22.5–23.0	1.0	0.0032 ± 0.0003	6.4	3.7	11
E0.7% $d3$				0.0037 ± 0.0002	11.1	5.5	7

$Re_{ul}$ —entrance Reynolds number,  $Ar_{ul}$ —Archimedes number in an upper layer,  $U_{ul}$ —terminal settling velocity, SD—standard deviation



Since fluid density, buoyancy frequency, and settling velocity are variable with depth for a particle translating in a nonlinear stratification, it is sensible to consider parameters as a function of depth. In such conditions, Reynolds number reads:

$$\text{Re}(z) = \frac{du(z)}{\nu}. \quad (3)$$

Archimedes number for a disk with characteristic length taken as equivalent sphere diameter (sphere of the same volume) is (Auguste et al. 2013):

$$\text{Ar}(z) = \frac{d}{\nu} \sqrt{\frac{3}{16} gh \frac{\rho_p - \rho_f(z)}{\rho_f(z)}} \quad (4)$$

where  $\rho_p$  is the particle density [ $\text{kg m}^{-3}$ ].

Froude number within transition layer reads:

$$\text{Fr}(z) = \frac{u(z)}{N(z)d}. \quad (5)$$

A ratio between Reynolds and Froude number forms another useful parameter (Mercier et al. 2020):

$$\text{Re}(z)/\text{Fr}(z) = \frac{N(z)d^2}{\nu}. \quad (6)$$

Settling velocity was estimated using a standard approach in homogeneous fluids to demonstrate to what extent oversimplified approach using formulas dedicated to homogeneous conditions may misestimate particle settling velocity and residence times in density transition. Settling velocity was evaluated iteratively using the equation for steady settling velocity of a disk:

$$u = \left( \frac{2gh}{C_d} \left( \frac{\rho_p}{\rho_f} - 1 \right) \right)^{0.5} \quad (7)$$

with drag coefficient,  $C_d$ , calculated from the formula dedicated to a terminal falling of a disk with Re number between 1.5 and 133 (Clift et al. 1978):

$$C_d = \frac{64}{\pi \text{Re}} (1 + 0.138 \text{Re}^{0.792}). \quad (8)$$

## Results and discussion

### Phases of disk settling and evolution of settling velocity

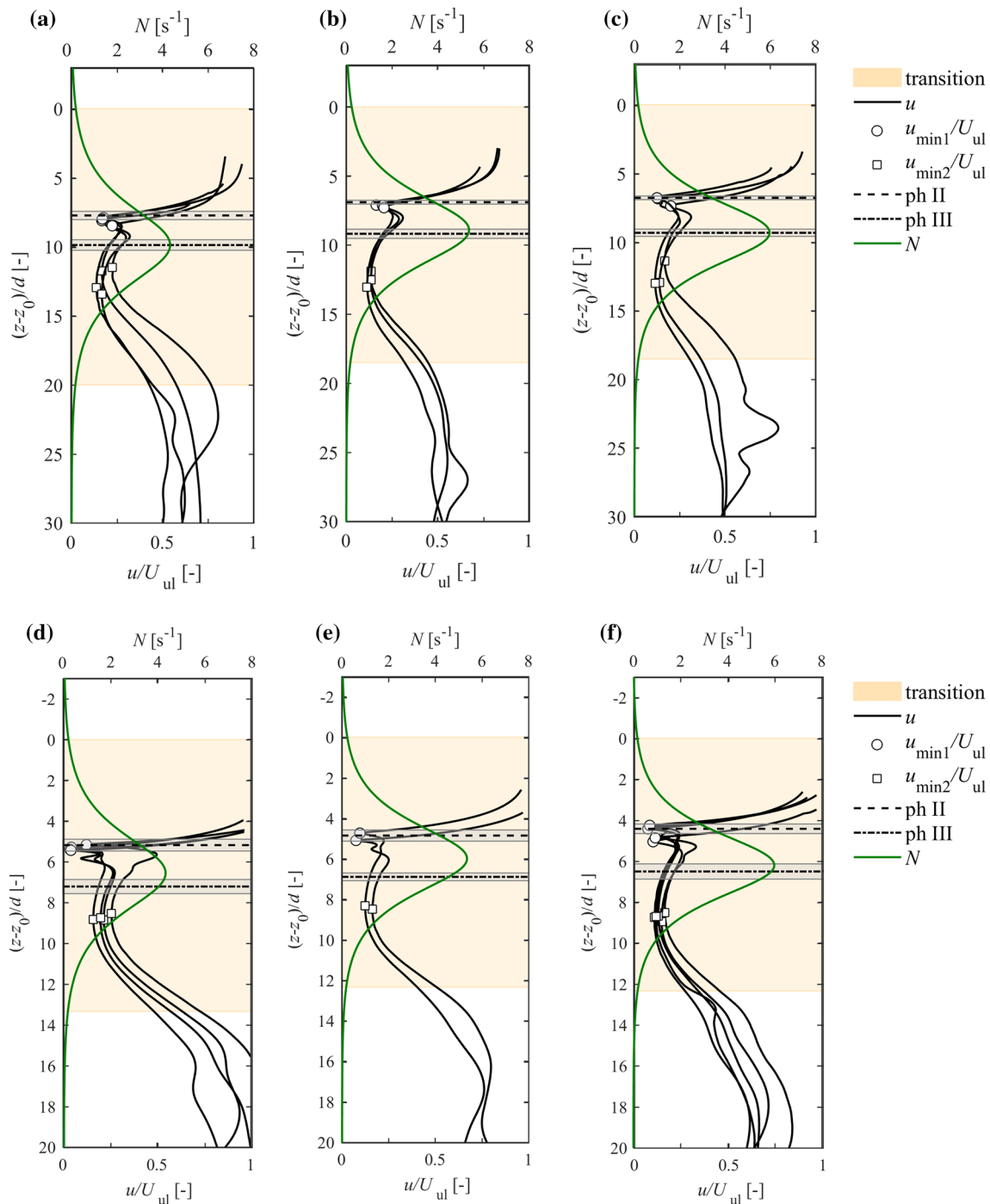
Disk settling dynamics followed the pattern observed in the previous research (Mrokowska 2018) where five phases of settling were identified (please refer to the explanation

therein). Figure 3 presents the typical variation of settling velocity with depth, where characteristic velocities, i.e. the first local velocity minimum,  $u_{\text{min}1}$ , the second local velocity minimum,  $u_{\text{min}2}$ , and mean locations for the beginning of phase II and III with corresponding standard deviations, are shown.

Disks were settling in an upper homogeneous layer (phase I) with terminal velocity varying from  $0.0025$  to  $0.0032 \text{ m s}^{-1}$  for  $d2$  disks and from  $0.0030$  to  $0.0037 \text{ m s}^{-1}$  for  $d3$  disks (not shown in the figure). Disks were translating broadside on, which is in line with phase diagrams describing dynamics of disks in terms of the relation between Re and dimensionless moment of inertia,  $I_*$  (Field et al. 1997; Willmarth et al. 1964) with the range of parameters observed in the study, Re between 5.0 and 11.1, and  $I_*$  of the order  $1 \times 10^{-3}$ , corresponding to the steady falling mode. Disks entered the transition with a broadside position and continued to settle in this orientation experiencing deceleration due to stratification effects.

Figure 3 shows that the location of the first minimum velocity,  $u_{\text{min}1}$ , corresponds with the beginning of reorientation (beginning of phase II). A particle assumes vertical position around the level where  $N = N_{\text{max}}$ , which may indicate that the maximum density gradient enhances rotation of particle to the vertical position. After the reorientation, the particle continues settling in a stable vertical position with a broadside perpendicular to horizontal (phase III) until stratification effects dominate over the inertia. It has been confirmed numerically in other study (Mercier et al. 2020) that stratification supports the descent of disk in a vertical position when settling velocity is low enough. Vertical position seems to be quite stable, and some particles descended in this position much further than others reaching the onset of reorientation to the horizontal position outside the camera field of view. It suggests that the orientation instability may be triggered not only by the stratification conditions but may be also induced by imperfections on the particle surface or in the location of the centre of mass, which may modify pressure distribution around the particle inducing rotation. Thus, “imperfect disks” are likely to change orientation earlier than perfect ones. The reorientation to the horizontal position occurs in a gliding motion with fading stratification effects when a disk is leaving the transition (phase IV). All studied disks achieved the second minimum velocity,  $u_{\text{min}2}$ , when translating in vertical position, which is analysed further in this paper.

Since phase IV (reorientation from vertical to horizontal position) was not recorded for some data sets due to limited field of view, data analysis for this phase is only partial. All particles assumed a stable horizontal position in a lower homogeneous layer (phase V); however, this phase was outside FOV and has not been presented herein.



**Fig. 3** Evolution of instantaneous settling velocity with depth for two-layered experiments **(a)** E2.6% $d_2$ , **(b)** E1.6% $d_2$ , **(c)** E0.7% $d_2$ , **(d)** E2.6% $d_3$ , **(e)** E1.6% $d_3$ , **(f)** E0.7% $d_3$  with the location of characteristic velocities  $u_{\min 1}$ ,  $u_{\min 2}$ . Vertical location of the onset of phases—the beginning of phase II (ph II) and the beginning of phase III (ph III)

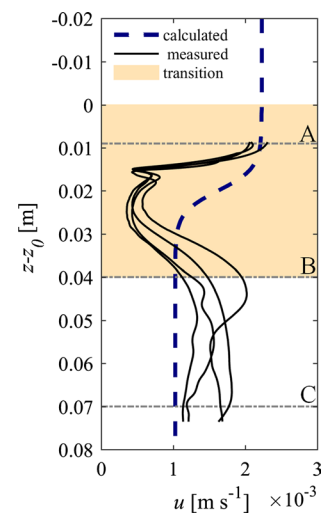
Data presented in Fig. 3 indicate that repeated experimental tests show a good agreement up to the point when particles achieve the level of maximum buoyancy frequency. This agreement is in both settling velocity values and the

location of reorientations. Below this level, the dispersion of settling velocity is observed and the onset of reorientation from vertical to horizontal position does not have a unique location.

All disks travelled some vertical distance in density transition before they started to reorient from a stable horizontal to vertical position. (See the location of the beginning of phase II in Fig. 3.) Similar effect has been observed in a linear stratification where it was found that stable broadside on position changes to vertical when a disk decelerates to some threshold velocity (Mercier et al. 2020). However, when we compare the settling velocity profile in a linear and nonlinear stratification, significant differences are evident. First of all, in a linear stratification velocity decreases monotonically (Mercier et al. 2020). Conversely, in the case of a nonlinear density gradient, a particle decelerates and accelerates in response to the variable stratification strength, which is additionally combined with the change of particle orientation. Consequently, drag exerted on a particle is a combination of buoyancy variation due to vertical density gradient, stratification strength, and axisymmetric particle shape. Quantification of drag is still challenging in such conditions and needs consideration in future studies.

### Effect of stratification on evaluation of settling velocity and residence time

Figure 4 presents the sample results of setting velocity estimation using Eq. (7) with drag coefficient defined with Eq. (8), the approach which is relevant in homogeneous conditions. The results are presented here to demonstrate that methods derived for a homogeneous fluid may result in serious misestimating of settling process in the presence of stratification. This approach showed satisfactory results in an upper homogeneous layer, slightly underestimating settling velocity with percentage error  $100(u_{\text{cal}} - u_{\text{obs}})/u_{\text{obs}}$ , where  $u_{\text{cal}}$  is the settling velocity calculated with Eq. (7) and  $u_{\text{obs}}$  is a measured value, between 4 and 10%. However, the comparison of the results obtained with Eq. (7) with velocities measured in the density-stratified region reveals significant deviations between estimated and observed values, which is shown in Fig. 4. First of all, this approach does not reproduce the existence of settling velocity minima. Consequently, predicted settling velocity in the major part of transition layer is larger than observed values. Secondly, this approach underestimates residence time of a particle in the transition layer (A–B in Fig. 4) as well as in the region where the particle was tracked (A–C region in Fig. 4). Results presented in Table 3 show that measured residence times in the transition are up to threefold larger than estimated by the homogeneous fluid approach. Moreover, the formula underestimates also the total residence time which varies between 0.61 and 0.88 of measured value.



**Fig. 4** Comparison between measured settling velocity and settling velocity calculated with Eq. (7) for sample data E2.6%*d*2. Grey dashed lines A, B, C show vertical locations described in text

### Trajectory and orientation pattern

Figure 5 shows the sample trajectories of disks retrieved from the view of settling particle recorded by camera#1 ( $x$ ,  $z$ ) and camera#2 ( $y$ ,  $z$ ) with the indication of characteristic velocity points and onsets of phases II, III, and IV. 2D positions of particle centre in the pairs of images, shown in Fig. 5b, c, e, f, h, i, were combined to get 3D trajectories of settling disks shown in corresponding plots in Fig. 5a, d, g. It could be seen from Fig. 5 that in phases I and III, when a disk settles in a horizontal and vertical position, respectively, a particle tends to settle in a vertical path. Some deviation from the perfect vertical path is observed as the effect of tilting of particles which could be due to unavoidable imperfections of manufactured particles (uneven surface, location of the centre of mass) which affect settling dynamics in a low-inertia motion. Nonetheless, settling could be considered as vertical. On the other hand, in reorientation phases II and IV, the particle moves in a horizontal plane as the effect of particle inclination with respect to the gravity.

The extent of particle horizontal drift was assessed as the Euclidean distance. Figure 6 presents the variation of distance travelled by the centre of particle mass with respect to its initial position in a horizontal plane. The results show very good repeatability of settling behaviour in phase II, while some discrepancies (deviation from vertical settling) are observed for phase III similarly to the above analysis of data presented in Fig. 5.

The analysis of 3D data shows that particles do not descend in a plane along the whole path (Figs. 5, 6). While the settling in phases I–III could be considered as planar, the disk may change the plane of settling in a second

**Table 3** Comparison between residence times of disks measured ( $t_{r, \text{obs}}$ ) and calculated ( $t_{r, \text{cal}}$ ) in the transition layer (A–B distance) and in A–C distance (see Fig. 4 for definition)

Exp	A–B distance			A–C distance		
	$t_{r, \text{cal}}$ [s]	$t_{r, \text{obs}}$ [s]	$t_{r, \text{cal}}/t_{r, \text{obs}}$ [–]	$t_{r, \text{cal}}$ [s]	$t_{r, \text{obs}}$ [s]	$t_{r, \text{cal}}/t_{r, \text{obs}}$ [–]
E2.6% <i>d</i> 2	24.1	42.9	0.56	55.3	63.0	0.88
E2.6% <i>d</i> 3	19.0	42.3	0.45	41.2	56.1	0.73
E1.6% <i>d</i> 2	22.3	48.0	0.46	55.7	72.7	0.77
E1.6% <i>d</i> 3	17.1	46.2	0.37	42.5	63.9	0.67
E0.7% <i>d</i> 2	20.3	40.6	0.50	51.7	62.4	0.83
E0.7% <i>d</i> 3	15.4	45.6	0.34	39.3	64.2	0.61

reorientation phase (phase IV), since there is no preferential direction of horizontal drift during this phase. It causes the translation of particle with respect to its initial position in a horizontal plane and consequently affects the extent of dispersion of particles.

Plots in Fig. 7 show the horizontal dispersion of particles in the  $x$ – $y$  plane. It is clear from the figure that there is no preferential direction of reorientations, which results in a random final position of particle with respect to the initial position. Horizontal drift is attributed to reorientation phases and has not exceeded 0.01 m along the vertical distance (0.077 m) analysed in this study. It should be noted that the fact that particles do not fall in a plane affects the variance of horizontal drift results.

The results reveal some trends in horizontal dispersion when particle diameter is considered. Figure 8 presents the comparison between an average horizontal drift for two sets of particles. The comparison between a drift in phases I and II and the total drift in considered paths indicates that the second reorientation contributes to the horizontal drift to a larger extent than the first one; that is, a particle travels larger horizontal distance in phase IV in which rotation is accompanied with gliding motion than in phase II. The results indicate that *d*3 disks tend to travel larger distance in a horizontal direction compared to *d*2 disks; hence, greater dispersion may be expected for disks with a larger diameter given the same stratification conditions. The effect of stratification strength on particle dispersion is not clear.

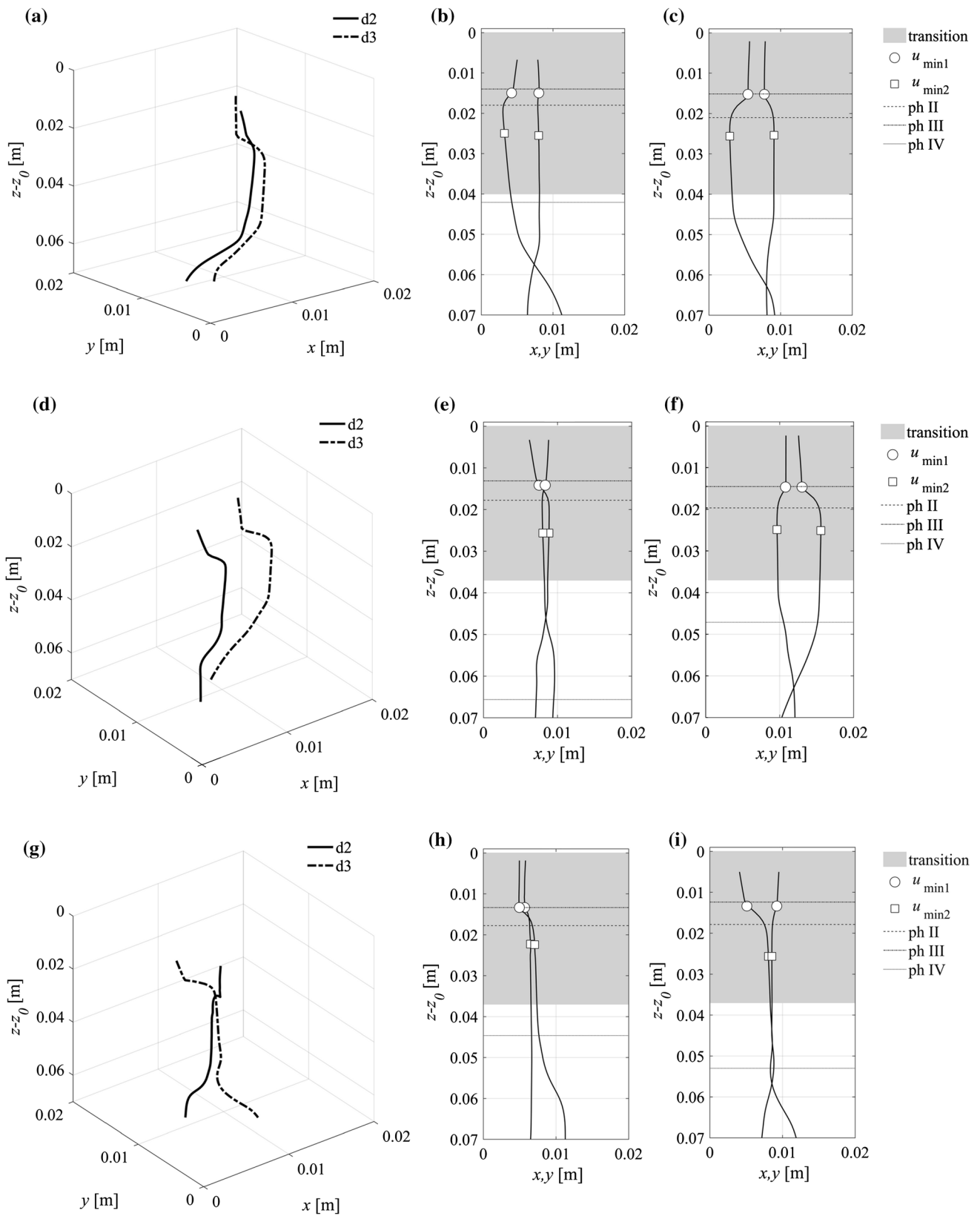
### Analysis of parameters for characteristic minimum velocities and onsets of reorientation

Settling velocity decreases significantly when a particle is translating across a density transition from terminal settling value in an upper layer,  $U_{ul}$ , to the first minimum settling velocity,  $u_{\text{min}1}$ , within the transition layer. Velocity reduction is more pronounced for larger particles (see Figs. 9 and 10) with  $u_{\text{min}1}/U_{ul}$  ranging between 0.05 and 0.08 for disks *d*3, while it is within the range (0.14; 0.19) for *d*2. Moreover, the ratio between  $u_{\text{min}1}$  for *d*3 and *d*2 is less than 75%, indicating that *d*3 achieves smaller local settling velocity after crossing the interface than *d*2. It could be explained

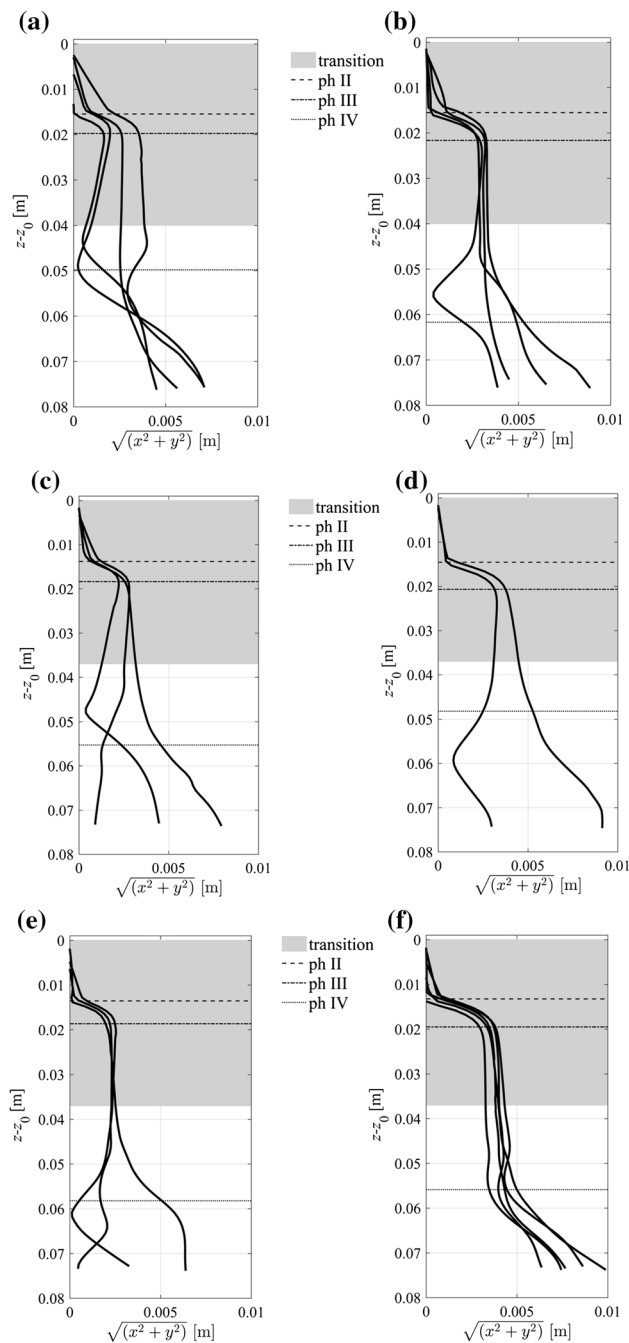
by the fact that a disk falling with its face normal to gravity in the upper part of transition region starts to experience additional stratification-induced drag. This drag is associated with the entrainment of lighter fluid into the wake of settling disk (see Fig. 3 in (Mrokowska 2018)). Since a disk with larger diameter is able to entrain larger volume of lighter fluid, the added-buoyancy effect increases with the particle dimensions, which implies larger drag and velocity reduction. On the other hand, local minimum velocity  $u_{\text{min}2}$  does not vary significantly between the two types of particles for the same conditions in density transition, indicating that settling velocity is not so sensitive on particle geometry when it descends in vertical position.

It is expected that the parameters for characteristic points within the density transition (the first and the second local minima and reorientation points) depend on the settling characteristics in the upper layer and the characteristics of stratification. Four depth-dependent parameters describing settling conditions within the transition layer relevant in this study problem have been analysed (Eqs. 3–6):  $Re(z)$ ,  $Fr(z)$ ,  $Ar(z)$ ,  $Re(z)/Fr(z)$  and the entrance Re number,  $Re_{ul}$  was considered to elucidate how settling dynamics is affected by the stratification and characteristics of settling dynamics in the upper layer.

To get insight into the particle dynamics, the evolution of Re number with depth has been analysed. While Re number assumes constant values in the upper layer between 5.0 and 11.1 depending on experimental conditions (particle diameter and fluid density) (Table 2), it varies with depth within the transition layer following the pattern of settling velocity. Re number achieves two local minima corresponding with settling velocity minima dropping to the minimum  $Re \sim 1$  for all considered particles. Reynolds number corresponding to the first minimum velocity,  $Re_{\text{um}in1}$ , achieves smaller values for larger disks;  $Re_{\text{um}in1}$  is within the range (0.92, 1.1) for *d*3 and within the range (0.48, 0.89) for *d*2 disks. Related observation for velocity is presented in Fig. 10 and described above. Observations for  $Re_{\text{um}in2}$  are in line with that for  $u_{\text{min}2}$ , that is, disk dimensions do not have significant impact on  $Re_{\text{um}in2}$ . Results presented in this study show that the onset of reorientation from the stable horizontal to stable vertical position overlaps with the minimum settling



**Fig. 5** Reconstructed 3D trajectories of disk descent and 2D  $(x,z)$  and  $(y,z)$  planar projections with indication of density transition, characteristic velocities points, and location of phases for sample experimental tests **(a, b, c)** E2.6%, **(d, e, f)** E1.6%, **(g, h, i)** E0.7%



**Fig. 6** Horizontal drift of disks in  $x$ - $y$  plane. (a) E2.6%,  $d_2$ , (b) E2.6%,  $d_3$ , (c) E1.6%,  $d_2$ , (d) E1.6%,  $d_3$ , (e) E0.7%,  $d_2$ , (f) E0.7%,  $d_3$

velocity,  $u_{\min 1}$ . Since Re number drops significantly at this point to values close to unity or even smaller, the reduction in inertial effects enables buoyancy-induced torques to overcome pressure-induced torques. Particles accelerate when descending in the transition in a vertical position as a result of diminishing stratification strength (Fig. 3), and Re number increases accordingly. However, disks are able to keep stable vertical position up to relatively high Re (in some tests  $\sim 6$ )

despite growing effect of inertia, which indicates that density gradients effectively overcome pressure-induced torques.

Another relevant parameter, Froude number, varies with depth between infinity in homogeneous layers to finite values within the transition layer. Fr achieves the minimum in the region where buoyancy frequency assumes the highest values and the particle decelerates considerably, indicating that the stratification dominates over inertia. Considering all experimental sets, minimum Froude number varies between 0.05 and 0.06 for  $d_2$  and between 0.02 and 0.03 for  $d_3$ , suggesting greater influence of stratification on a disk with a larger diameter. Fr number increases up to  $Fr \sim 5$  at the lower border of transition where  $N = 0.2 \text{ s}^{-1}$ , and then it increases to infinity in a homogeneous layer.

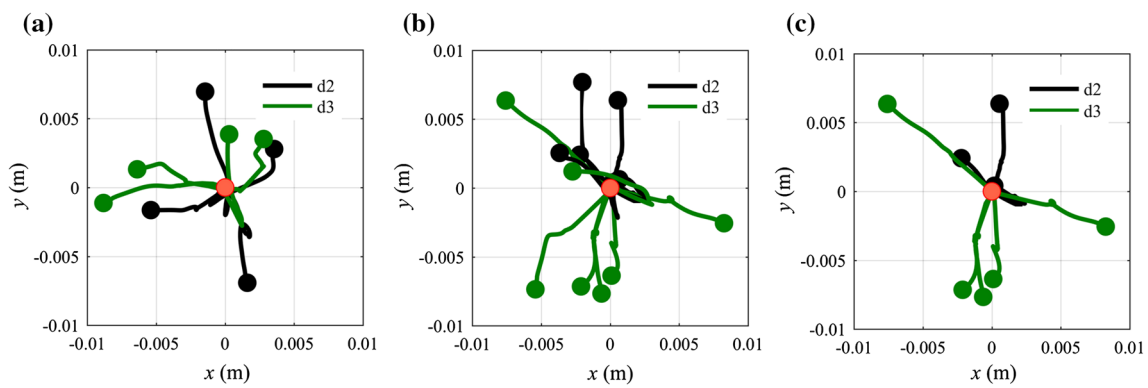
Archimedes number decreases monotonically from a constant value in an upper layer (between 3.5 and 6.2 in considered experimental sets) to 2.2 and 3.3 for  $d_2$  and  $d_3$ , respectively, in the lower layer showing a great dependence on density difference between the fluid and the particle.

The relation between settling dynamics in an upper layer, characterised by  $Re_{ul}$ , and characteristic values of parameters  $Fr(z)$ ,  $Re(z)/Fr(z)$ ,  $Ar(z)$  corresponding to settling velocity minima, denoted here by a subscript “c”, were analysed, and the results are presented in Fig. 11. Plots show a decreasing trend for  $Fr_{\min 1}$  versus  $Re_{ul}$  and increasing trend for  $Ar_{\min 1}$  versus  $Re_{ul}$  and  $Re_{\min 1}/Fr_{\min 1}$  versus  $Re_{ul}$ , indicating that these relations may be good candidates to elucidate threshold values for the onset of the first reorientation (and occurrence of the first minimum velocity). These results indicate that settling dynamics in the upper layer, stratification characteristics in the transition layer, and particle dimensions affect characteristic minimum velocities and the onset of reorientation which is compatible with  $u_{\min 1}$ . However, more data sets for wider range of conditions are necessary to define a functional relationship.

Similar relationships are observed for parameters  $Fr_{\min 2}$ ,  $Re_{\min 2}/Fr_{\min 2}$ , and  $Ar_{\min 2}$ ; however, trends are weaker, showing that the occurrence of the second velocity minimum is less sensitive to the settling conditions in the upper layer than the first minimum velocity. No clear dependence was observed between characteristic Reynolds number and entrance Re number.

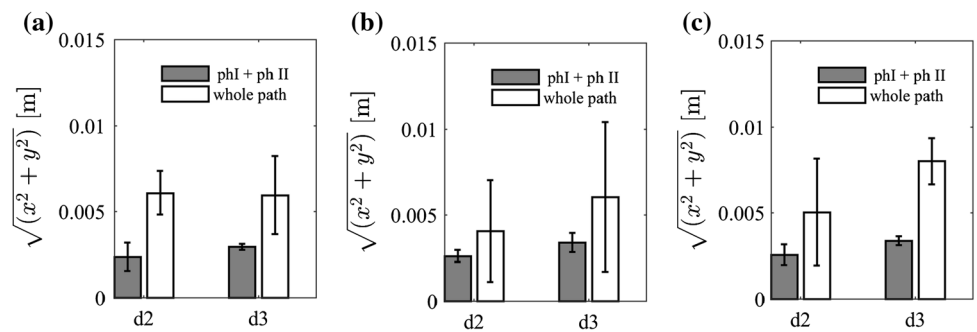
## Conclusions

Settling dynamics of thin disks descending through a non-linear density transition were studied experimentally. The results have demonstrated that complex hydrodynamic interactions between a particle and a liquid lead to settling orientation instabilities and unsteady particle descent. The most interesting aspects of disk settling through the density transition are conditions necessary for reorientation,

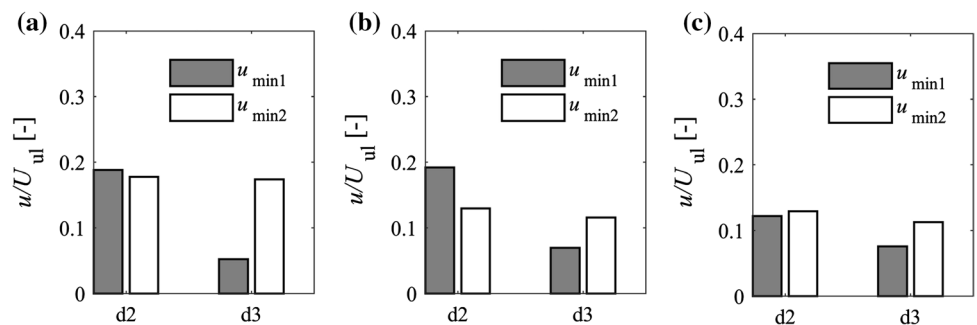


**Fig. 7** Projection of disk trajectory in  $x$ - $y$  plane showing planar dispersion of particles. Initial position of disk centre of mass marked by a red dot has been set as  $(0, 0)$  and final recorded points marked by green and black dots for disks  $d3$ , and  $d2$ , respectively

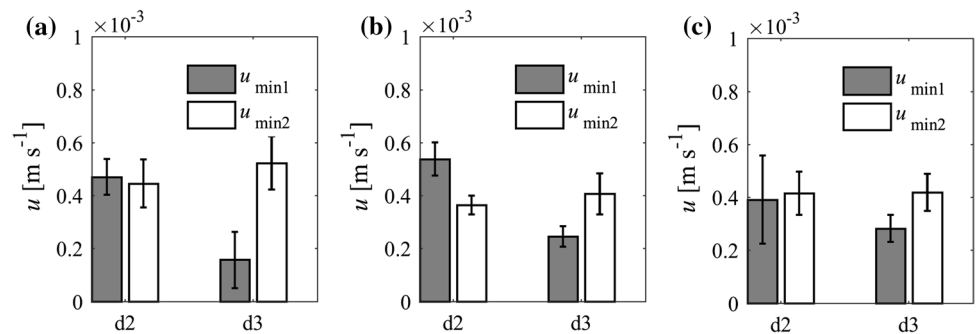
**Fig. 8** Average horizontal drift in phase I and II and in the whole path for disks  $d2$  and  $d3$ ; **(a)** E2.6%, **(b)** E1.6%, **(c)** E0.7%

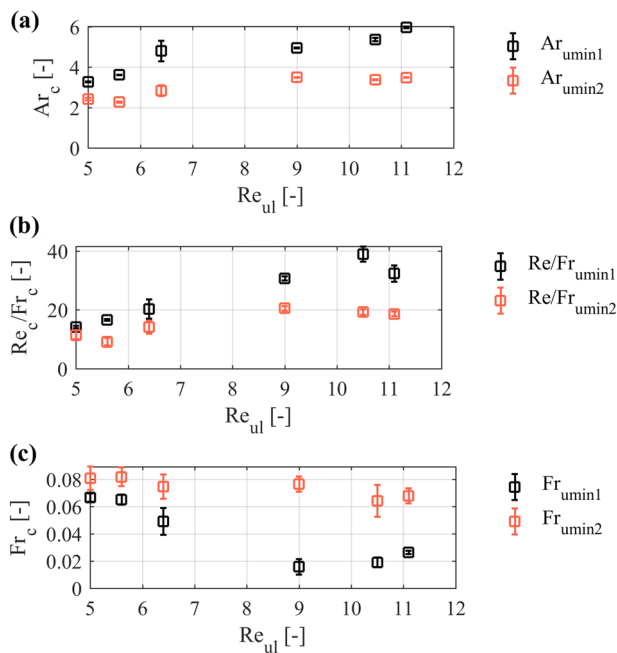


**Fig. 9** Local velocity minima  $u_{\min 1}$  and  $u_{\min 2}$  in relation to terminal settling velocity in an upper layer,  $U_{ul}$ , for disks  $d2$  and  $d3$  in two-layered experiments **(a)** E2.6%, **(b)** E1.6%, **(c)** E0.7%



**Fig. 10** Average values for local velocity minima  $u_{\min 1}$  and  $u_{\min 2}$  for disks  $d2$  and  $d3$  in two-layered experiments, error bars denote standard deviation **(a)** E2.6%, **(b)** E1.6%, **(c)** E0.7%





**Fig. 11** Relation between settling conditions in an upper layer characterised by entrance  $Re$  number,  $Re_{ul}$ , and critical values of parameters for minimum velocities. Mean results for all experiments are presented with standard error

the pattern of non-monotonous velocity with focus on local minima in the second and third phase, and the horizontal drift of particles.

Relations between entrance Reynolds number,  $Re_{ul}$ , and  $Re/Fr$ ,  $Ar$ , and  $Fr$  were identified for characteristic points (settling velocity local minima and the first reorientation points), indicating settling conditions in an upper layer, the stratification strength of transition, and particle dimension control settling dynamics within the transition layer.

The condition for reorientation from a stable horizontal to stable vertical position is of particular interest, which has been analysed for linear stratification in other research (Mercier et al. 2020). Results presented herein show that a particle keeps a horizontal position within nonlinear density transition up to a point where the particle achieves the first local minimum and the onset of reorientation has been attributed to  $Re$  number low enough to let stratification effects overcome inertial ones.

The comparison of settling dynamics of thin disks with two different diameters showed that the dynamics are sensitive to a particle diameter mainly in the upper part of density transition with a non-obvious result that a larger disk settling with higher terminal velocity in the upper layer achieves smaller first minimum velocity than disk with a smaller diameter. Settling dynamics of two types of disks within a transition layer does not vary considerably in terms of settling velocity with negligible difference in the second

minimum velocity. All these results suggest that geometry of particles should be carefully considered to assess settling dynamics just after crossing a density interface, since small deviations in geometry may affect settling parameters.

Fundamental understanding of dynamics of individual particles in stratified systems is critical for further elucidation of physical mechanisms of particle groups settling necessary to develop prediction methods for sedimentation flux and descent of immotile microorganisms. Specifically, the knowledge on settling dynamics within transition layer may improve the understanding of thin layers formation, since one of reasons for little understanding of this process is our scarce knowledge on the dynamics of non-spherical particles and effects of stratification. The results showed that disks experience horizontal drift in reorientation phases (phase II and IV) and fading stratification in phase IV and a gliding motion of particle in this phase play the major role in horizontal dispersion of particles. This fact may be important in a wider context when a group of disks settles and may interact with each other, e.g. a group of disk-shaped diatoms in the ocean. Since the dispersion of particles is expected in the lower part of transition, interactions between particles are likely to be intensified in this region.

It has been demonstrated that parameters critical to estimate particulate flux, namely the residence time of particles in a water column and settling velocity, are misestimated by conventionally used approaches which do not take into account the dynamics of non-spherical particles in stratified systems. Settling behaviour of particles within the transition layer has impact not only on the prolonged residence of particles at pycnoclines which may lead to the formation of thin layers, but also increases the total residence time of particles in a water column, which is of importance to the estimation of sedimentation rate, carbon transport, as well as other biogeochemical processes.

Density transitions occurring commonly in natural waters are much thicker compared to the dimensions of settling particles and could be considered as linear; nonetheless, the results presented herein reveal physical mechanisms that could explain settling behaviour of solid particles in nature. Although a vast amount of particles present in the environment is of non-spherical shapes including disks, a few research performed so far (Doostmohammadi and Ardekani 2014; Mercier et al. 2020; Mrokowska 2018) indicated that the settling dynamics including orientation instabilities and non-monotonous settling velocity are more important than probably previously thought. Hence, present challenge is to incorporate the geometry of non-spherical particles into the analysis of settling behaviour of particles in background stratification, since the variation of particle orientation with its descent highly affects settling velocity due to orientation-dependent drag.



This study is constrained to the limited experimental conditions, i.e. particle characteristics and stratification parameters; only a disk diameter was variable while thickness was kept constant. The relationships between settling parameters could be different for disks with different aspect ratio which should be considered in further studies. More experimental results on disks settling through a density transition are necessary to find parameter ranges characteristic for various behaviours of disks settling in a low to moderate Re numbers. Wider range of experimental conditions should be considered to investigate some limiting cases in density transition with a nonlinear stratification to quantify conditions governing the onset of stratification-induced reorientations. Moreover, development of numerical models is necessary to facilitate laboratory experiments due to their limited capacity.

**Author contributions** Magdalena Mrokowska is the sole author of the present study; she conceived the study, performed experiments, processed data, analysed results, and wrote the paper.

**Funding** This work was supported within an internal grant of Institute of Geophysics, Polish Academy of Sciences No. 7a/IGF PAN/2018ml and partially within statutory activities No. 3841/E-41/S/2020 of the Ministry of Science and Higher Education of Poland.

**Data availability** Data available upon request to the author.

## Compliance with ethical standards

**Conflict of interest** The author declares no competing interests.

**Code availability** Not applicable.

**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

## References

- Abaid N, Adalsteinsson D, Agyapong A, McLaughlin RM (2004) An internal splash: levitation of falling spheres in stratified fluids. *Phys Fluids* 16:1567–1580. <https://doi.org/10.1063/1.1687685>
- Ardekani AM, Doostmohammadi A, Desai N (2017) Transport of particles, drops, and small organisms in density stratified fluids. *Phys Rev Fluids* 2:100503. <https://doi.org/10.1103/PhysRevFluids.2.100503>
- Arnosti C (2011) Microbial extracellular enzymes and the marine carbon cycle. *Annu Rev Mar Sci* 3:401–425. <https://doi.org/10.1146/annurev-marine-120709-142731>
- Auguste F, Magnaudet J, Fabre D (2013) Falling styles of disks. *J Fluid Mech* 719:388–405. <https://doi.org/10.1017/jfm.2012.602>
- Bagheri G, Bonadonna C (2016) On the drag of freely falling non-spherical particles. *Powder Technol* 301:526–544. <https://doi.org/10.1016/j.powtec.2016.06.015>
- Blanchette F, Shapiro AM (2012) Drops settling in sharp stratification with and without Marangoni effects. *Phys Fluids* 24:042104. <https://doi.org/10.1063/1.4704790>
- Camassa R, Falcon C, Lin J, McLaughlin RM, Mykins N (2010) A first-principle predictive theory for a sphere falling through sharply stratified fluid at low Reynolds number. *J Fluid Mech* 664:436–465. <https://doi.org/10.1017/s0022112010003800>
- Camassa R, Falcon C, Lin J, McLaughlin RM, Parker R (2009) Prolonged residence times for particles settling through stratified miscible fluids in the Stokes regime. *Phys Fluids* 21:031702. <https://doi.org/10.1063/1.3094922>
- Capet A, Stanev EV, Beckers JM, Murray JW, Gregoire M (2016) Decline of the Black Sea oxygen inventory. *Biogeosciences* 13:1287–1297. <https://doi.org/10.5194/bg-13-1287-2016>
- Clift R, Grace JR, Weber ME (1978) Bubbles, drops, and particles. Academic Press, Cambridge
- Cole M, Lindeque P, Halsband C, Galloway TS (2011) Microplastics as contaminants in the marine environment: a review. *Mar Pollut Bull* 62:2588–2597. <https://doi.org/10.1016/j.marpolbul.2011.09.025>
- D'Asaro E (2018) Oceanographic floats: principles of operation. In: Venkatesan R, Tandon A, D'Asaro E, Atmanand MA (eds) *Observing the oceans in real time*. Springer, Cham, pp 77–98. [https://doi.org/10.1007/978-3-319-66493-4\\_5](https://doi.org/10.1007/978-3-319-66493-4_5)
- Dey S, Ali SZ, Padhi E (2019) Terminal fall velocity: the legacy of Stokes from the perspective of fluvial hydraulics. *Proc R Soc A-Math Phys Eng Sci* 475:20190277. <https://doi.org/10.1098/rspa.2019.0277>
- Diercks A, Ziervogel K, Sibert R, Joye SB, Asper V, Montoya JP (2019) Vertical marine snow distribution in the stratified, hypersaline, and anoxic Orca Basin (Gulf of Mexico). *Elementa-Sci Anthropol* 7:1. <https://doi.org/10.1525/elementa.348>
- Dietrich WE (1982) Settling velocity of natural particles. *Water Resour Res* 18:1615–1626. <https://doi.org/10.1029/WR018i006p01615>
- Doostmohammadi A, Ardekani AM (2014) Reorientation of elongated particles at density interfaces. *Phys Rev E* 90:033013. <https://doi.org/10.1103/PhysRevE.90.033013>
- Doostmohammadi A, Dabiri S, Ardekani AM (2014) A numerical study of the dynamics of a particle settling at moderate Reynolds numbers in a linearly stratified fluid. *J Fluid Mech* 750:5–32. <https://doi.org/10.1017/jfm.2014.243>
- Doostmohammadi A, Stocker R, Ardekani AM (2012) Low-Reynolds-number swimming at pycnoclines. *Proc Natl Acad Sci USA* 109:3856–3861. <https://doi.org/10.1073/pnas.1116210109>
- Field SB, Klaus M, Moore MG, Nori F (1997) Chaotic dynamics of falling disks. *Nature* 388:252–254. <https://doi.org/10.1038/40817>
- Kestin J, Khalifa HE, Correia RJ (1981) Tables of the dynamic and kinematic viscosity of aqueous NaCl solutions in the temperature-range 20–150-degrees-C and the pressure range 0.1–35 MPa. *J Phys Chem Ref Data* 10:71–87
- Khatmullina L, Isachenko I (2017) Settling velocity of microplastic particles of regular shapes. *Mar Pollut Bull* 114:871–880. <https://doi.org/10.1016/j.marpolbul.2016.11.024>
- Kindler K, Khalili A, Stocker R (2010) Diffusion-limited retention of porous particles at density interfaces. *Proc Natl Acad Sci USA* 107:22163–22168. <https://doi.org/10.1073/pnas.1012319108>

- Lam T, Vincent L, Kanso E (2019) Passive flight in density-stratified fluids. *J Fluid Mech* 860:200–223. <https://doi.org/10.1017/jfm.2018.862>
- Laurenceau-Cornec EC, Le Moigne FAC, Gallinari M et al (2019) New guidelines for the application of Stokes' models to the sinking velocity of marine aggregates. *Limnol Oceanogr* 9999:1–22. <https://doi.org/10.1002/lno.11388>
- Loth E (2008) Drag of non-spherical solid particles of regular and irregular shape. *Powder Technol* 182:342–353. <https://doi.org/10.1016/j.powtec.2007.06.001>
- Lutz M, Dunbar R, Caldeira K (2002) Regional variability in the vertical flux of particulate organic carbon in the ocean interior. *Global Biogeochem Cycles* 16(3):1037. <https://doi.org/10.1029/2000gb001383>
- Macintyre S, Alldredge AL, Gotschalk CC (1995) Accumulation of marine snow at density discontinuities in the water column. *Limnol Oceanogr* 40:449–468. <https://doi.org/10.4319/lno.1995.40.3.0449>
- Maggi F (2013) The settling velocity of mineral, biomineral, and biological particles and aggregates in water. *J Geophys Res-Oceans* 118:2118–2132. <https://doi.org/10.1002/jgrc.20086>
- Magnaudet J, Mercier MJ (2020) Particles, drops, and bubbles moving across sharp interfaces and stratified layers. *Annu Rev Fluid Mech* 52:61–91. <https://doi.org/10.1146/annurev-fluid-010719-060139>
- Mercier MJ, Wang S, Pemeja J, Ern P, Ardekani AM (2020) Settling disks in a linearly stratified fluid. *J Fluid Mech* 885:A2. <https://doi.org/10.1017/jfm.2019.957>
- Mrokowska MM (2018) Stratification-induced reorientation of disk settling through ambient density transition. *Sci Rep* 8:412. <https://doi.org/10.1038/s41598-017-18654-7>
- Mrokowska MM, Krztoń-Maziopa A (2019) Viscoelastic and shear-thinning effects of aqueous copolymer solution on disk and sphere settling. *Sci Rep* 9:7897. <https://doi.org/10.1038/s41598-019-44233-z>
- Noufal KK, Najeem S, Latha G, Venkatesan R (2017) Seasonal and long term evolution of oceanographic conditions based on year-around observation in Kongsfjorden, Arctic Ocean. *Polar Sci* 11:1–10. <https://doi.org/10.1016/j.polar.2016.11.001>
- Okino S, Akiyama S, Hanazaki H (2017) Velocity distribution around a sphere descending in a linearly stratified fluid. *J Fluid Mech* 826:759–780. <https://doi.org/10.1017/jfm.2017.474>
- Peperzak L, Colijn F, Koeman R, Gieskes WWC, Joordens JCA (2003) Phytoplankton sinking rates in the Rhine region of freshwater influence. *J Plankton Res* 25:365–383. <https://doi.org/10.1093/plankt/25.4.365>
- Prairie JC, White BL (2017) A model for thin layer formation by delayed particle settling at sharp density gradients. *Cont Shelf Res* 133:37–46. <https://doi.org/10.1016/j.csr.2016.12.007>
- Prairie JC, Ziervogel K, Camassa R et al (2015) Delayed settling of marine snow: effects of density gradient and particle properties and implications for carbon cycling. *Mar Chem* 175:28–38. <https://doi.org/10.1016/j.marchem.2015.04.006>
- Prairie JC, Ziervogel K, Camassa R et al (2017) Ephemeral aggregate layers in the water column leave lasting footprints in the carbon cycle. *Limnol Oceanogr Lett* 2:202–209. <https://doi.org/10.1002/lo12.10053>
- Raffaele L, Bruno L, Sherman DJ (2020) Statistical characterization of sedimentation velocity of natural particles. *Aeol Res*. <https://doi.org/10.1016/j.aeolia.2020.100593>
- Renggli CJ, Wiesmaier S, De Campos CP, Hess KU, Dingwell DB (2016) Magma mixing induced by particle settling. *Contrib Miner Petrol* 171(11):96. <https://doi.org/10.1007/s00410-016-1305-1>
- Saxby J, Beckett F, Cashman K, Rust A, Tennant E (2018) The impact of particle shape on fall velocity: Implications for volcanic ash dispersion modelling. *J Volcanol Geoth Res* 362:32–48. <https://doi.org/10.1016/j.jvolgeores.2018.08.006>
- Scase MM, Dalziel SB (2004) Internal wave fields and drag generated by a translating body in a stratified fluid. *J Fluid Mech* 498:289–313. <https://doi.org/10.1017/s0022112003006815>
- Srdic-Mitrovic AN, Mohamed NA, Fernando HJS (1999) Gravitational settling of particles through density interfaces. *J Fluid Mech* 381:175–198. <https://doi.org/10.1017/s0022112098003590>
- Turner JT (2015) Zooplankton fecal pellets, marine snow, phytodetritus and the ocean's biological pump. *Prog Oceanogr* 130:205–248. <https://doi.org/10.1016/j.pocean.2014.08.005>
- Verso L, van Reeuwijk M, Liberzon A (2019) Transient stratification force on particles crossing a density interface. *Int J Multiph Flow*. <https://doi.org/10.1016/j.ijmultiphaseflow.2019.103109>
- Waldschlager K, Schuttrumpf H (2019) Effects of particle properties on the settling and rise velocities of microplastics in freshwater under laboratory conditions. *Environ Sci Technol* 53:1958–1966. <https://doi.org/10.1021/acs.est.8b06794>
- Willmarth WW, Hawk NE, Harvey RL (1964) Steady and unsteady motions and wakes of freely falling disks. *Phys Fluids* 7:197–208. <https://doi.org/10.1063/1.1711133>
- Woods AW (1995) The dynamics of explosive volcanic-eruptions. *Rev Geophys* 33:495–530. <https://doi.org/10.1029/95rg02096>
- Yick KY, Torres CR, Peacock T, Stocker R (2009) Enhanced drag of a sphere settling in a stratified fluid at small Reynolds numbers. *J Fluid Mech* 632:49–68. <https://doi.org/10.1017/s0022112009007332>
- Zhai L, Sun ZB, Li ZM et al (2019) Dynamic effects of topography on dust particles in the Beijing region of China. *Atmos Environ* 213:413–423. <https://doi.org/10.1016/j.atmosenv.2019.06.029>



# Trend analysis and SARIMA forecasting of mean maximum and mean minimum monthly temperature for the state of Kerala, India

P. Kabbilawsh<sup>1</sup> · D. Sathish Kumar<sup>1</sup> · N. R. Chithra<sup>1</sup>

Received: 5 February 2020 / Accepted: 3 July 2020 / Published online: 16 July 2020  
© Institute of Geophysics, Polish Academy of Sciences & Polish Academy of Sciences 2020

## Abstract

The development of temperature forecasting models for the state of Kerala using Seasonal Autoregressive Integrated Moving Average (SARIMA) method is presented in this article. Mean maximum and mean minimum monthly temperature data, for a period of 47 years, from seven stations, are studied and applied to develop the model. It is expected that the time-series datasets of temperature to display seasonality (and hence non-stationary), and a possible trend (due to the fact that the data spans 5 decades). Hence, the key step in the development of the models is the determination of the non-stationarity of the temperature time-series, and the transformation of the non-stationary time-series into a stationary time-series. This is carried out using the Seasonal and Trend decomposition using Loess technique and Kwiatkowski–Phillips–Schmidt–Shin test. Before carrying out this process, several preliminary tests are conducted for (1) finding and filling the missing values, (2) studying the characteristics of the data, and (3) investigating the presence of the trend and seasonality. The non-stationary temperature time-series are transformed to stationary temperature time-series, by one seasonal differencing and one first-order differencing. This information, along with the original time-series, is further utilized to develop the models using the SARIMA method. The parsimonious and best-fit SARIMA models are developed for each of the fourteen variables. The study revealed that SARIMA(2, 1, 1)(1, 1, 1)<sub>12</sub> as the ideal forecasting model for eight out of the fourteen time-series datasets.

**Keywords** Autocorrelation function (ACF) · Partial autocorrelation function (PACF) · Sen's slope estimator · Seasonal autoregressive integrated moving average (SARIMA) · Mann–Kendall (MK) trend test

## Introduction

India, with a population of more than 1.3 billion, has more than 50% of its population dependent on agriculture (Arjun 2013). Most states in India still heavily rely on rainfall for various agricultural activities. It is well known that rainfall, a part of the hydrological cycle, is susceptible to changes in global temperature (Allen and Ingram 2002; Andronova and Schlesinger 2000; Trenberth 1999). Hence, an exclusive look into the long-term temperature variations would

constitute a vital part in the analysis of agricultural output of any region of the country.

In this regard, many researchers have carried out studies in the last decade on global, continental and regional level long-term temperature variations (Hänsel et al. 2016; Jain and Kumar 2012; Kocsis et al. 2017). Also, many attempts have been undertaken by researchers to develop models for understanding and extrapolating the temperature variation (Hänsel et al. 2016; Mills 2014; Tiwari et al. 2016). In India, among all the studies focused on temporal temperature variation, the most noteworthy study is the one conducted by the Indian Network for Climate Change Assessment (INCCA) (2010). The projections of mean annual surface temperature for the 2030s (average of 2021–2050) were carried out on country level using PRECIS (Providing Regional Climates for Impact Studies), with the data obtained from 1970s (average of 1960–1990). In this study, it was predicted that the annual mean surface air temperature would rise by 1.7–2 °C over the entire Indian subcontinent. Though this study indicates that significant

✉ P. Kabbilawsh  
kabbi.civil@gmail.com

D. Sathish Kumar  
sathish@nitc.ac.in

N. R. Chithra  
chithranr@nitc.ac.in

<sup>1</sup> Department of Civil Engineering, NIT Calicut, Calicut 673601, India

changes could be expected in the overall characteristics of rainfall, the projections are at a macroscopic level (i.e. for the entire Indian subcontinent), and not for each individual states. Regional studies focussing on individual states are necessary to get a better understanding of the local factors that influence these variations. A state-wise study is important because local policies and actions can be exclusively implemented by the state governments to combat any expected adverse changes in their respective states. In this study, the temporal variation of the monthly mean maximum (MMAX) and mean minimum temperature (MMIN) is analysed for the state of Kerala.

The analysis is carried out for a period of 47 years, starting from 1969 to 2015. The overall objective of the study is to develop a model for future forecasts of MMAX and MMIN for the state of Kerala. Prior to the time-series modelling, it is necessary to carry out preprocessing of the data to identify the missing values. The time-series data available for each station and the number of missing values are listed in Table 1. The data gaps are to be eliminated before any time-series modelling. The data infilling process is carried out using expectation–maximization algorithm. Further, for the construction of a suitable forecasting model, it is necessary to evaluate the time-series datasets to understand the existing pattern. This preliminary analysis provides a good insight regarding the available data. It comprises of (1) a descriptive statistical analysis of the monthly data, (2) performance of the normality test, (3) test to check for outliers, (4) Mann–Kendall trend analysis and (5) performance of the Sen’s slope test. The results obtained from the preliminary analysis revealed the presence of non-stationarity in the datasets. In order to confirm the preliminary results obtained, the time-series datasets are decomposed using STL decomposition to get the time-series components. The obtained time-series components also revealed the presence of seasonality and the presence or absence of a trend. The value of parameters (seasonal and non-seasonal differencing,  $D$  and  $d$ , respectively) needed for converting the non-stationary time-series to a stationary series is obtained using the results of Kwiatkowski–Phillips–Schmidt–Shin (KPSS) test. These values, along with the original time-series datasets, are used for the SARIMA model building process. The next section (“Temperature data and research methodology” section) briefly describes each of the process (the preliminary tests, STL decomposition, Unit root test and SARIMA) applied in this study. Section “Temperature data and research methodology” describes the application of these tests to our data. Also, in “Result and discussions” section, the result of each test is analyzed and elaborated, and a final forecast is also delivered with the developed model. Lastly, section “Summary and

conclusions” concludes the article with an overview of the entire study.

## Temperature data and research methodology

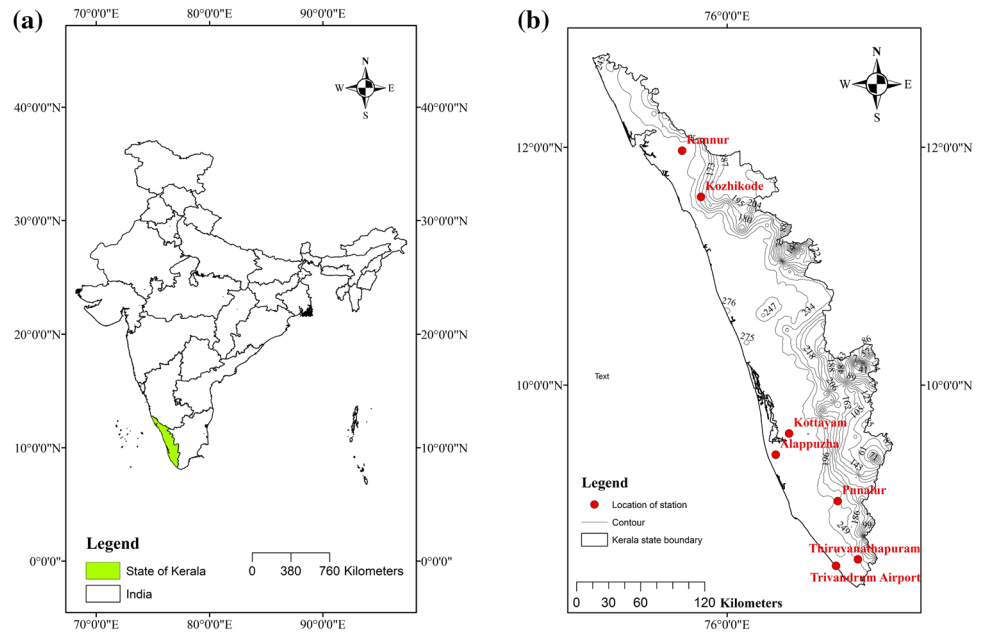
The main reason for carrying out the temperature-related studies for the state of Kerala is that the state is the gateway of the summer monsoon (South-West) for India. Any disturbance to the South-West monsoon creates a cascading effect on the rainfall patterns in the entire country. As stated earlier, this section presents the techniques applied for (1) estimating the missing values, (2) conducting the preliminary analysis, (3) decomposing the time-series data, (4) converting the non-stationary data to stationary data, and (5) developing the model.

### An account of Kerala and its temperature dataset

The state Kerala is a small strip of coastal land located in the southern part of India. It consists of an area of 38,850 km<sup>2</sup>. It is located between 8° 18′ N–12° 48′ N latitudes and 74° 52′ E–77° 24′ E longitudes. Figure 1 shows the location map of the study area. The state has a shoreline of 580 km, and the width of the state varies between 30 and 120 km. Geologically, the state Kerala can be categorised into three climatically distinct regions: the eastern highlands (rugged and cool mountainous terrain), the central midlands (rolling hills), and the western lowlands (coastal plains). The lowlands and highlands bound the state of Kerala, where the lowlands comprise the regions which adjoin the shoreline, and highlands cover the region slopping down from the Western Ghats. The midlands spread between the highlands and lowlands. Area-wise, the highlands comprise of 18,650 km<sup>2</sup>, while the midlands and lowlands comprise of 16,200 km<sup>2</sup> and 4000 km<sup>2</sup> respectively. Tea, coffee and rubber are major plantation crops grown in the highlands. It also houses several endemic flora and fauna. Wide variety of fruits, nuts and vegetables are grown in the midland region. Paddy and coconut are grown in the fertile lowlands.

The temperature datasets from 13 observatories that cover the entire state of Kerala are obtained from the Indian Meteorological Department (IMD). The observatories spread across all the three climatic regions. The observatories are located at Palghat, Fort Cochin, Kovalam, Karipur, Trichur, Ernakulum, Kozhikode, Kannur, Alappuzha, Punalur, Kottayam, Thiruvananthapuram and Trivandrum Airport. Out of 13 stations, three observatories (Palghat, Fort Cochin and Kovalam) are not properly functioning for the past 15 years, and the datasets for recent years are not available. For three observatories (Karipur, Trichur and Ernakulum), the datasets are available starting from the year 1996 and

**Fig. 1** **a** Location map of the state Kerala, **b** the location map of the seven stations for which study is conducted



**Table 1** The amount of missing data present in the meteorological observatories

Station name	Starting year of time series	Ending year of time series	Total length of data (years)	MMAX			MMIN		
				Number of monthly values present	Number of monthly values missing	% of missing values	Number of monthly values present	Number of monthly values missing	% of missing values
Kozhikode	1969	2015	47	564	0	0	564	0	0
Kannur	1981	2015	34	404	16	3.81	404	16	3.81
Alappuzha	1969	2015	47	549	15	2.66	548	16	2.84
Punalur	1969	2015	47	532	32	5.67	500	64	11.34
Kottayam	1973	2011	43	498	18	3.49	496	20	3.88
Thiruvananthapuram	1969	2015	47	564	0	0	564	0	0
Trivandrum Airport	1969	2015	47	545	19	3.37	545	19	3.37

later. Adequate datasets for the analysis are available only for seven meteorological stations. The spatial locations of the observatories are shown in Fig. 1b.

The data from the rest of the seven stations are found ideal for the study. Table 1 shows the data availability for the selected seven stations. A total of 235 intermittent monthly values (about 3.13% of the data) are found to be missing in the available dataset. Table 1 lists the number and types of missing values for each station. Datasets with missing values present several problems in the representativeness of the samples (Kang 2013). Hence, the missing values are to be determined first. For this purpose, the expectation–maximization algorithm is used. The missing values estimated through this method is used to fill the data gaps in order to obtain continuous time-series datasets. Preliminary

statistical tests are conducted using these datasets. The results indicated the presence of skewness and kurtosis. Further, a test for normality is carried out using the Shapiro-Wilk normality test and the outliers are identified using Grubb’s test. The results indicated that datasets followed a non-normal distribution without any outliers. Therefore, a nonparametric Mann–Kendall trend test (Gocic and Trajkovic 2013; Kocsis et al. 2020) and Sen’s slope test are used to determine the direction and magnitude of monotonic trends in the time-series.

### Mann–Kendall trend test

The Mann–Kendall trend test (Mann 1945; Kendall 1975; Gilbert 1987) is widely used test in the field of

Hydro-meteorology, dealing with variables like temperature, rainfall and streamflow. The Mann–Kendall test is used to statistically assess the presence of an increasing or decreasing trend in the series. The Mann–Kendall test operates by checking whether to reject the null hypothesis ( $H_0$ ) and accept the alternative hypothesis ( $H_1$ ). The null hypothesis ( $H_0$ ) means there is no trend in the temperature over time, and the alternative hypothesis ( $H_1$ ) implies the presence of either an increasing or decreasing trend in the temperature data. The sign of the computed Mann–Kendall test statistic  $Z_{MK}$  reveals the direction of the trend. The positive value of  $Z_{MK}$  indicates that the temperature tends to increase with time, while the negative value of  $Z_{MK}$  denotes the decrease in temperature over time. The null Hypothesis ( $H_0$ ) is rejected, and the alternative hypothesis ( $H_1$ ) is accepted if  $|Z_{MK}| \geq Z_{1-\alpha/2}$  at the Type I error rate  $\alpha$ .

### Sen's slope estimation

All the available statistical techniques may not be equally good in detecting the magnitude of the trend in the time-series data (Radziejewski and Kundzewicz 2004). A simple parametric least-square regression technique is not suitable to calculate the magnitude of the trend for non-normal time-series. In such cases, a test which is nonparametric, robust against outliers would be an appropriate choice. Sen's Slope estimation, a nonparametric test is selected to detect the magnitude of trends in the temperature time-series. It is impartially resistant to outliers, with a breakdown point of 0.29 (Sen 1968). It was initially proposed in 1968 to account for the non-normality of precipitation data. A mathematical explanation of the scheme is not detailed here, as it has been already presented in detail by various authors (Gocic and Trajkovic 2013; Kocsis et al. 2020).

### STL decomposition

The Mann–Kendall trend test and Sen's slope estimation are carried out as a part of the initial investigation. Further to provide a better understanding of the datasets, the time-series data is decomposed as the trend component, the seasonal component and the remainder component. It is carried out using STL (Seasonal and Trend decomposition using Loess) decomposition method (Cleveland et al. 1990). The decomposed components are plotted for graphical visualisation of the data. It allows us to visualise the presence of trend and seasonality in the data. Compared to the other classical decomposition methods, STL has several advantages like the ability to handle any type of seasonality (daily, monthly, quarterly, annual, etc.), being robust to outliers, facilitating the user to control the smoothness of trend cycle, and allowing the user to control the rate of change of seasonal component.

### Unit root test

The trend and seasonal components obtained from the STL decomposition will reveal the presence of non-stationarity in the temperature time-series. The non-stationarity is only inferred from the graphs of the decomposed components (only visual inference). To mathematically confirm the presence of non-stationarity in the time-series, the unit root tests are performed. In the present study, KPSS (Kwiatkowski–Phillips–Schmid–Shin) unit root test is performed to confirm the presence of non-stationarity (Kwiatkowski et al. 1992). The original temperature time-series and the decomposed components are used for this purpose. The KPSS method proceeds with the null hypothesis (i.e. the data are stationary) and tries to find evidence to show that the null hypothesis is false for the selected time-series. If the non-stationarity is confirmed, then the next step is the conversion of the non-stationary data to stationary data. The  $p$  values determined from the KPSS test provides information about the differencing; small  $p$  values typically, less than 0.05 points the necessity of differencing for the conversion of the time-series.

### Seasonal autoregressive integrated moving average (SARIMA) model

After the non-stationary time-series is converted to a stationary time-series (i.e. after the determination of  $d$  and  $D$ ), the next step is to develop a model for future predictions. Forecasting models developed from the historical records are generally used to predict the future changes in the climate variables. Several authors have proposed temperature models using a number of forecasting techniques (Aguado-Rodríguez et al. 2016; Tiwari et al. 2016; Wang et al. 2019; Lai and Dzombak 2020; Wanishsakpong and Owusu 2020). Several climate variables are generally influenced by seasonality, and one of the best forecasting models for such variables is the SARIMA model. It combines the advantage of the autoregressive model and the moving average model.

In an autoregressive model, a linear combination of the past values of the variable is used to predict the future of the variable. Mathematically, Eq. 1 represents an autoregressive model of order  $p$ , i.e.  $AR(p)$  model.

$$y_t = \theta_0 + \theta_1 y_{t-1} + \theta_2 y_{t-2} + \dots + \theta_p y_{t-p} + \varepsilon_t \quad (1)$$

The equation shows that the observation  $y$  at time  $t$  ( $y_t$ ) is estimated from  $p$  previous observations ( $y_{t-i}$ ,  $i = 1, 2, 3, \dots, p$ ).  $\theta_k$ , with  $k = 1, 2, 3, \dots, p$  are the parameters, and  $\varepsilon_t$  is the white noise.

In a moving average model, the forecast is done using the past forecast errors in a regression-like model.

$$y_t = \phi_0 + \varepsilon_t + \phi_1\varepsilon_{t-1} + \phi_2\varepsilon_{t-2} + \dots + \phi_q\varepsilon_{t-q} \tag{2}$$

Equation 2 describes the moving average model of order  $q$ , i.e.  $MA(q)$  model, where  $y_t$  is the observation at time  $t$ ;  $\phi_i$ , with  $i = 1, 2, 3, \dots, q$ , are the parameters and  $\varepsilon_{t-k}$ , with  $k = 1, 2, 3, \dots, q$  are the error terms, respectively.

By combining differencing with autoregression and a moving average model, the non-seasonal Autoregressive Integrated Moving Average (ARIMA) model is obtained. Mathematically, the ARIMA model is represented by Eq. 3.

$$y'_t = \theta_0 + \theta_1y'_{t-1} + \theta_2y'_{t-2} + \dots + \theta_p y'_{t-p} + \phi_1\varepsilon_{t-1} + \phi_2\varepsilon_{t-2} + \dots + \phi_q\varepsilon_{t-q} + \varepsilon_t \tag{3}$$

$y'_t$  denotes the differenced series. It has to be noted that the series may have been differenced more than once, and the degree of the differencing involved is denoted by  $d$ . This series is represented as the  $ARIMA(p, d, q)$  model, where  $p$  denotes the order of autoregressive part,  $d$  denotes the degree of differencing, and  $q$  denotes the moving average part.

If seasonality is observed in a time-series, then a seasonal-ARIMA model or SARIMA model (Hyndman and Athanasopoulos 2018) has to be applied. The seasonal-ARIMA model is obtained by including the additional seasonal terms to the ARIMA models. The seasonal-ARIMA model is represented as  $SARIMA(p, d, q)(P, D, Q)_m$ . The non-seasonal part of the model is represented as  $(p, d, q)$ , and the seasonal part of the model is given by  $(P, D, Q)_m$ . The terms  $P, D$  and  $Q$  represents the order of the seasonal autoregressive term, degree of the seasonal differencing and order of the seasonal moving average part, respectively. The term  $m$  represents the number of observations per year. The terms of the seasonal part of the model are similar to the

non-seasonal part expect that they involve backshifts of the seasonal period.

## Result and discussions

All the tests discussed in the previous section are applied parallelly or sequentially based on the requirements. The results of these tests, and their significance, are discussed in this section.

### Results from descriptive statistics

Descriptive statistics of the temperature datasets are obtained after filling the missing values using the expectation–maximization algorithm. The analysis is carried out for seven stations for two variables (MMAX and MMIN) in each station. Therefore, altogether fourteen time-series datasets are analysed. The average MMAX and MMIN temperature covering all the seven stations are 31.84 °C and 23.48 °C respectively. The MMAX varies between 31.11 °C (at Trivandrum Airport), and 33.06 °C (at Punalur), and MMIN varies between 22.34 °C (at Punalur) and 24.22 °C (at Kozhikode). The standard error of the mean of all fourteen variables ranges between 0.04 °C and 0.1 °C. The maximum deviation of the sample-mean from the population-mean is 0.2 °C, at a confidence level of 95%. The difference between the sample-mean and the population-mean is negligible. Therefore, it can be concluded that a sample mean is a genuine representation of the population mean. The descriptive statistics of the variables are listed in Table 2. In the tabulation, SEM,

**Table 2** Descriptive statistics of the variables

Station name	Variable type	Mean	SEM	SD	Variance	CV	$Q_1$	$Q_3$	Range	IQR	Skewness	Excess kurtosis
Kozhikode	MMAX	31.50	0.08	1.79	3.22	5.69	30.20	32.80	8.40	2.60	−0.23	−0.64
Kozhikode	MMIN	24.22	0.05	1.24	1.54	5.12	23.50	24.80	8.10	1.30	0.22	0.29
Kannur	MMAX	32.13	0.10	2.04	4.16	6.35	30.50	33.70	8.90	3.20	−0.14	−0.89
Kannur	MMIN	23.47	0.07	1.33	1.78	5.68	22.70	24.20	7.20	1.50	0.12	−0.12
Alappuzha	MMAX	31.48	0.07	1.64	2.70	5.22	30.10	32.80	7.70	2.70	−0.25	−0.90
Alappuzha	MMIN	23.92	0.05	1.17	1.36	4.88	23.20	24.60	6.50	1.40	0.07	0.03
Punalur	MMAX	33.06	0.09	2.14	4.58	6.47	31.40	34.70	10.30	3.30	0.36	−0.57
Punalur	MMIN	22.34	0.05	1.19	1.42	5.33	21.70	23.10	7.10	1.40	−0.25	0.14
Kottayam	MMAX	32.03	0.08	1.76	3.11	5.51	30.70	33.40	9.50	2.70	−0.09	−0.65
Kottayam	MMIN	23.08	0.04	0.98	0.95	4.23	22.70	23.70	6.20	1.00	−0.91	1.52
Thiruvananthapuram	MMAX	31.56	0.06	1.38	1.90	4.37	30.50	32.68	6.20	2.18	0.04	−0.83
Thiruvananthapuram	MMIN	23.56	0.04	0.97	0.94	4.12	23.00	24.10	5.90	1.10	0.24	0.13
Trivandrum Airport	MMAX	31.11	0.05	1.18	1.39	3.79	30.20	31.90	5.90	1.70	0.21	−0.66
Trivandrum Airport	MMIN	23.76	0.05	1.10	1.20	4.61	23.20	24.32	6.60	1.12	−0.23	0.60

SD, CV,  $Q$  and IQR stand for standard error of the mean, standard deviation, coefficient of variation, quartile, and inter-quartile range, respectively.

The MMIN variable at Kottayam has a skewness of  $-0.91$ , and since this value is within the range of  $-0.5$  and  $-1$ , it implies that data is moderately skewed. Moreover, the excess kurtosis values of all fourteen variables are nonzero, which implies that all fourteen temperature time-series are non-mesokurtic. Although the excess kurtosis values (nonzero) indicate the non-normal distribution of all the variables, it has to be noted that the values are small. Therefore, it necessitates a dedicated normality test. Consequently, a Shapiro-Wilk test is conducted to validate the nature of the distribution.

## Test for normality

The test for normality indicated that all fourteen variables are indeed non-normally distributed. The results of the Shapiro-Wilk test are presented in Table 3.

Additionally, the Grubb's test is also conducted to determine the presence of outliers in the data. The results of the Grubb's test are presented in Table 4. The G-statistic values of all fourteen variables are found to be less than their corresponding critical values indicating that there are no outliers in any of the fourteen temperature datasets.

**Table 3** The results of the test for normality of the variables

Station name	Variable type	Degrees of freedom	Shapiro–Wilk		
			Statistic	$p$ value	Decision at level (5%)
Kozhikode	MMAX	564	0.983	3E–06	Reject normality
Kozhikode	MMIN	564	0.977	9E–08	Reject normality
Kannur	MMAX	420	0.978	6E–06	Reject normality
Kannur	MMIN	420	0.987	1E–03	Reject normality
Alappuzha	MMAX	564	0.972	6E–09	Reject normality
Alappuzha	MMIN	564	0.991	2E–03	Reject normality
Punalur	MMAX	564	0.979	3E–07	Reject normality
Punalur	MMIN	564	0.989	4E–04	Reject normality
Kottayam	MMAX	516	0.988	2E–04	Reject normality
Kottayam	MMIN	516	0.953	1E–11	Reject normality
Thiruvananthapuram	MMAX	564	0.984	7E–06	Reject normality
Thiruvananthapuram	MMIN	564	0.983	5E–06	Reject normality
Trivandrum Airport	MMAX	564	0.983	4E–06	Reject normality
Trivandrum Airport	MMIN	564	0.981	1E–06	Reject normality

**Table 4** The results obtained from the Grubb's test for the variables

Station name	Variable type	G-statistic	Critical value	Approximate $p$ value (%)	Decision
Kozhikode	MMAX	2.4	3.9	9.18	No outliers
Kozhikode	MMIN	3.49	3.9	0.26	No outliers
Kannur	MMAX	2.22	3.82	10.87	No outliers
Kannur	MMIN	2.95	3.82	1.3	No outliers
Alappuzha	MMAX	2.61	3.9	5.05	No outliers
Alappuzha	MMIN	2.93	3.9	1.85	No outliers
Punalur	MMAX	2.82	3.9	2.61	No outliers
Punalur	MMIN	1.83	2.89	1.78	No outliers
Kottayam	MMAX	2.97	3.87	1.49	No outliers
Kottayam	MMIN	3.78	3.87	0.01	No outliers
Thiruvananthapuram	MMAX	2.35	3.9	10.4	No outliers
Thiruvananthapuram	MMIN	3.26	3.9	0.6	No outliers
Trivandrum Airport	MMAX	2.88	3.9	2.22	No outliers
Trivandrum Airport	MMIN	3.34	3.9	0.45	No outliers



## Trend analysis using Mann–Kendall test and Sen’s slope estimation

Since the datasets are not normally distributed, the Mann–Kendall test is applied to check the presence or absence of the trend in the datasets. The results of the trend analysis are presented in Table 5. As mentioned earlier, the results of the Mann–Kendall indicate only the presence or the absence of a trend in the series and its direction. However, it fails to quantify the magnitude of the trend.

In this test, a  $p$  value greater than  $\alpha$  (i.e. 0.05), indicates the absence of the trend. The sign of MK-statistic indicates the direction of the trend. The test results indicate that a certain amount of trend is present in ten variables. The magnitude of trend determined using Sen’s slope estimation is presented in Table 5. The  $\beta$ -slope represents the magnitude of the trend. This is consistent with the findings ( $p$  value) from the Mann–Kendall’s test. The four stations that indicated the absence of a trend in the Mann–Kendall test resulted in very low values of  $\beta$ -slope. It may be noted that for non-stationary series with small slopes ( $< 0.0002$ ), even at  $p < 0.01$ , Mann–Kendall trend test rejects null-hypothesis, resulting in Type-I error. The other specific inferences that could be made from this test is that the MMAX series of Kozhikode, Kannur and Thiruvananthapuram has a significant trend ( $\beta$ -slope exceeding 0.2%), and MMIN series of Alappuzha station is the only one with a decreasing trend, confirming the result obtained from the Mann–Kendall test.

## Analysis through STL decomposition

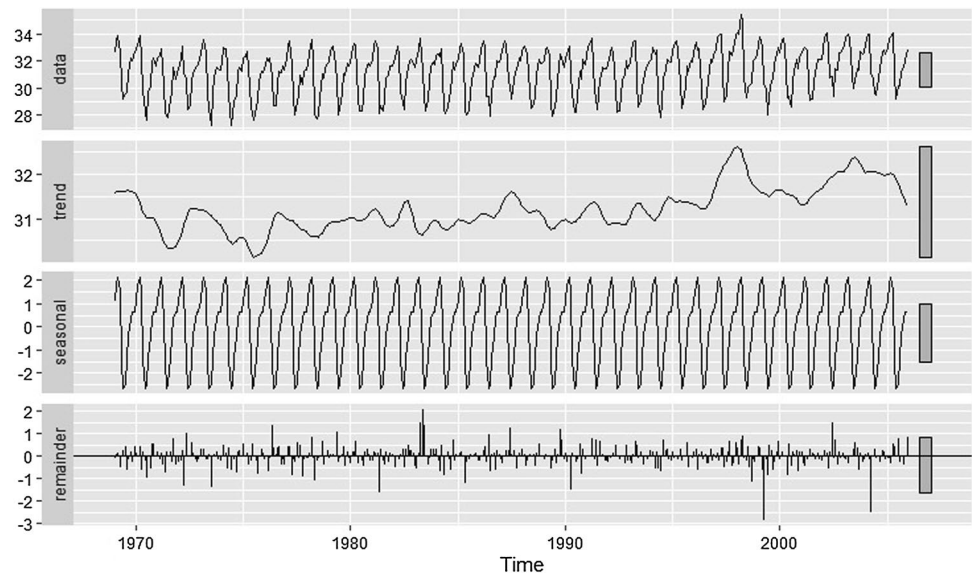
Though the previous two tests reveal the presence or absence of a trend, the presence of non-stationarity resulting from seasonality cannot be directly inferred from them. A time-series decomposition technique is applied to obtain the components. The fourteen time-series variables are decomposed to get each one’s trend, seasonal and remainder components using the STL decomposition. Figure 2 shows the original data, trend, seasonal, and remainder components of the MMAX variable of Kozhikode station.

For STL decomposition, only the first 80% of the time-series datasets are utilized. The remaining 20% of the data is retained for the validation of the forecasting model. The increasing trend, which was predicted by both the Mann–Kendall test and Sen’s slope estimation, can be visualised from the figure. It also indicates the existence of a strong seasonal pattern. Similar to Kozhikode MMAX variable, the other thirteen variables also exhibited seasonal patterns. On this basis, it is possible to conclude that the datasets are non-stationary. Before developing a forecasting model, the non-stationary datasets must be converted to stationary datasets, and subsequently the parameters  $d$  and  $D$  must be determined. The non-stationarity of the datasets are validated by applying the unit root test.

**Table 5** The results of Mann–Kendall trend test ( $\alpha = 0.05$ ) and Sen’s slope test ( $\alpha = 0.05$ )

Station name	Variable type	MK S-statistic	Standard error	$z$ statistic	$p$ value	Presence of trend	Sen’s slope	Sen’s-slope (lower 95 % Confidence Interval)	Sen’s-slope (Upper 95 % Confidence Interval)
Kozhikode	MMAX	35,826	4469.78	8.02	0	Yes	0.0038	0.0029	0.0046
Kozhikode	MMIN	26,011	4467.88	5.82	0	Yes	0.0015	0.001	0.0021
Kannur	MMAX	16,583	2873.8	5.77	0	Yes	0.0048	0.0032	0.0063
Kannur	MMIN	9118	2873	3.17	0	Yes	0.0016	0.0006	0.0025
Alappuzha	MMAX	9169	4469.78	2.05	0.04	Yes	0.0009	0	0.0017
Alappuzha	MMIN	– 14,724	4468.65	– 3.3	0	Yes	– 0.001	– 0.0016	– 0.0004
Punalur	MMAX	2670	4470.12	0.6	0.55	No	0.0003	– 0.0008	0.0015
Punalur	MMIN	5220	4468.91	1.17	0.24	No	0.0003	– 0.0001	0.001
Kottayam	MMAX	5403	3912.08	1.38	0.17	No	0.0007	– 0.0003	0.0018
Kottayam	MMIN	– 4099	3909.71	– 1.05	0.3	No	0	– 0.0008	0
Thiruvananthapuram	MMAX	33,767	4469.45	7.56	0	Yes	0.0028	0.0021	0.0035
Thiruvananthapuram	MMIN	20,561	4467.24	4.6	0	Yes	0.001	0.0006	0.0015
Trivandrum airport	MMAX	18,572	4469.01	4.16	0	Yes	0.0013	0.0007	0.0019
Trivandrum airport	MMIN	21,046	4467.87	4.71	0	Yes	0.0011	0.0006	0.0016

**Fig. 2** The decomposed components of the Kozhikode MMAX time-series



**Unit root test and the conversion to a stationary series**

The Kwiatkowski–Phillips–Schmid–Shin (KPSS) test is applied to categorise the datasets as stationary or non-stationary. The results of the unit root test are presented in Table 6. The test results for the original temperature time-series are listed in the third column of the table. A *p* value of less than 0.05 implies that the series is non-stationary. The results indicate that time series datasets of Alappuzha, Punalur and Kottayam corresponding to MMAX variable, and MMIN variable of Punalur and Kottayam are stationary. This is contrary to what was inferred from the STL seasonal plots. To resolve this paradox, the autocorrelation (ACF)

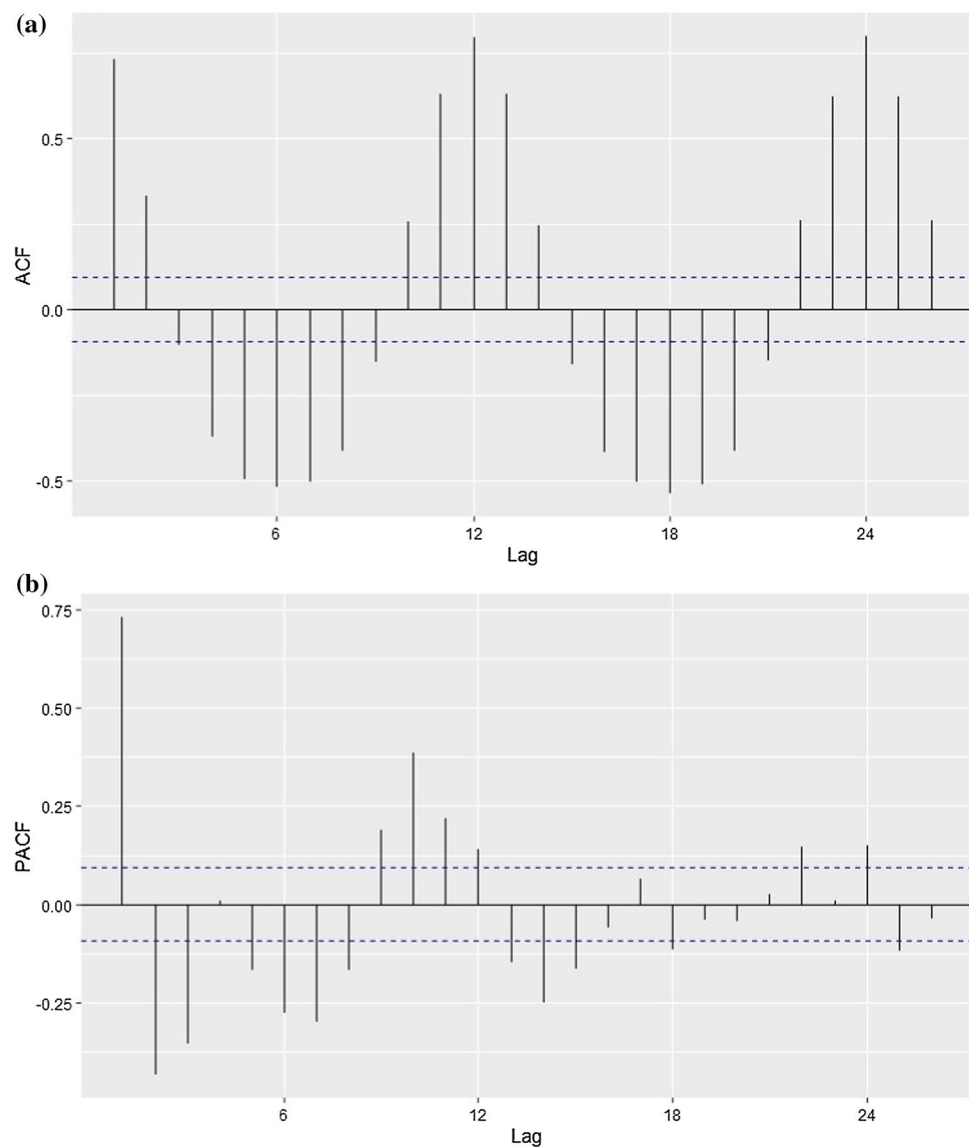
plots and partial autocorrelation (PACF) plots of those five variables are analyzed. The ACF and PACF plots of the Punalur MMAX are shown in Fig. 3. The ACF and PACF values in the plots follow a decaying sinusoidal pattern, which indicates seasonality and the dataset is non-stationary.

Similar trends are observed in the ACF and PACF plots for the other four variables. Therefore, it is conclusive that all fourteen variables are indeed non-stationary. As the seasonality is confirmed, at least one seasonal differencing is necessary to convert the non-stationary time-series to a stationary time-series. As this technique ascertains non-stationarity by means of seasonality, there are possibilities that the non-stationarity may exist exclusive of the seasonal component. In other words, it is possible that there could

**Table 6** The result of KPSS test for level stationarity

Station name	Variable type	Original series ( <i>p</i> value)	Seasonally adjusted series ( <i>p</i> value)	Series after seasonal adjustment and first-order difference ( <i>p</i> value)
Kozhikode	MMAX	0.01	0.01	0.1
Kozhikode	MMIN	0.01	0.01	0.1
Kannur	MMAX	0.01	0.01	0.1
Kannur	MMIN	0.049	0.01	0.1
Alappuzha	MMAX	0.1	0.01	0.1
Alappuzha	MMIN	0.01	0.01	0.1
Punalur	MMAX	0.1	0.1	0.1
Punalur	MMIN	0.1	0.051	0.1
Kottayam	MMAX	0.1	0.01	0.1
Kottayam	MMIN	0.06	0.01	0.1
Thiruvananthapuram	MMAX	0.01	0.01	0.1
Thiruvananthapuram	MMIN	0.01	0.01	0.1
Trivandrum Airport	MMAX	0.01	0.01	0.1
Trivandrum Airport	MMIN	0.01	0.01	0.1

**Fig. 3** **a** ACF plot of MMAX variable at Punalur station, **b** PACF plot of MMAX variable at Punalur station



be a non-stationarity in the non-seasonal component of the time-series. In order to determine this possibility, the KPSS test is conducted on the seasonally adjusted time-series.

The seasonally adjusted series is obtained by subtracting the seasonal component from the original time-series datasets (seasonal differencing). The test results for the seasonally adjusted time-series are listed in Table 6. The results indicate that most of the seasonally adjusted series are non-stationary. Hence, it is evident that the non-stationarity of the datasets is not just due to the presence of seasonality alone. It indicates that, in addition to seasonal differencing, performing the first-order difference would be prudent in conversion of non-stationary time-series to stationary time-series. Subsequently, all fourteen time-series datasets are subjected to one seasonal differencing and a first-order difference. The KPSS test is performed on the resulting time-series datasets. The test results are presented in the

last column of Table 6, where it can be observed that all the differenced time-series datasets are stationary.

### Modelling by seasonal autoregressive integrated moving average (SARIMA) method

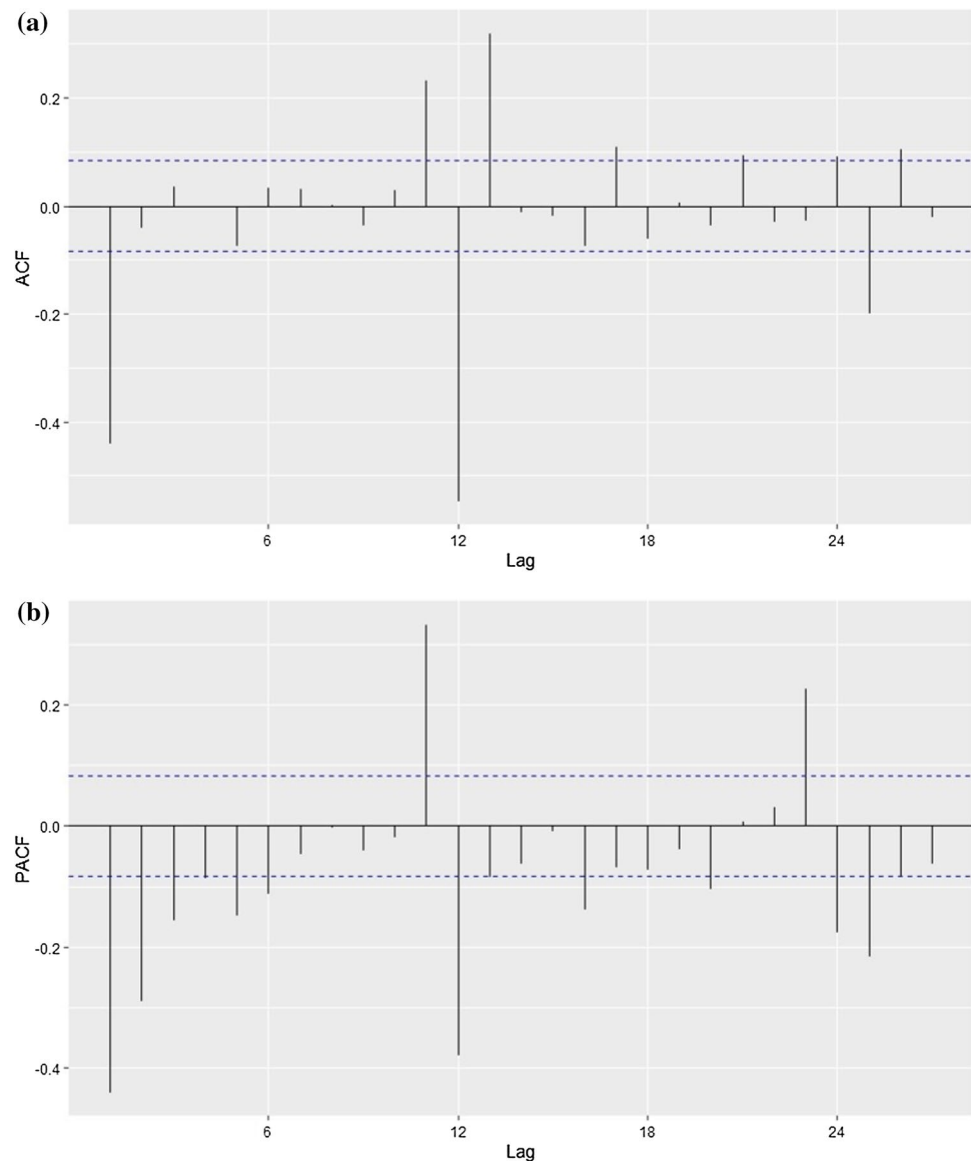
Forecasting models are developed by applying the SARIMA method using the original time-series datasets. The model has an inherent ability to transform non-stationary data into stationary data using the parameters  $d$  and  $D$  determined earlier. SARIMA models for each variable are developed with a different combination of parameters, and the best-fitting model is selected based on statistical evaluation. The procedure followed to develop the SARIMA model for MMAX variable of Kozhikode station is detailed. A similar procedure is adopted for the other thirteen variables.

In the SARIMA model, where the model is represented by  $SARIMA(p, d, q)(P, D, Q)_m$ , the determination of the parameters  $p, d, q, P, D, Q$  and  $m$ , for a particular time-series data, completes the development of the model for that dataset. In the present case, from the previous analysis, it was found that  $d = 1$  and  $D = 1$ , and for monthly data  $m = 12$ . Therefore, it is necessary to determine the values for parameters  $p, q, P$  and  $Q$  alone. These parameters are determined from the ACF and PACF plots of the stationary series (the seasonally and first-order differenced series) (Hyndman and Athanasopoulos 2018). The ACF and PACF plots of Kozhikode MMAX which are seasonally and first-order differenced are shown in Fig. 4, where the first 30 lags are considered for determining the parameters.

The value of non-seasonal autoregressive term ( $p$ ) and seasonal autoregressive term ( $P$ ) are determined from the

PACF plot (Fig. 4b). In the first span of seasonality, there are significant spikes at lag 1, lag 2 and lag 3, and this indicates that a non-seasonal autoregressive component up to  $AR(3)$  (i.e.  $p \leq 3$ ) would be appropriate. The spikes at lags 1, 2 and 3 are considered, while the spikes at lags 5, 6 and 11 are ignored, because lags 1, 2 and 3 serially lie outside the bounds and lag 4 lies within the bound, and thus break the continuity. All the out of bound lags, in the first span of seasonality, after lag 4 are ignored for this reason. Both the second (lags 12 to 23) and third span (lags 24 to 35) of seasonality have out of bound lags. Therefore, a seasonal autoregressive component  $AR(2)$  (i.e.  $P \leq 2$ ) would be appropriate. Similarly, the moving average components are determined from the ACF plot Fig. 4a. The appropriate values of moving average components are  $q \leq 1$  and  $Q \leq 2$  (seasonal). Thus, the candidate model is  $SARIMA(3, 1, 1)(2, 1, 2)_{12}$ .

**Fig. 4** **a** ACF plots for the stationary series of MMAX variable at Kozhikode station, **b** PACF plot for the stationary series of MMAX variable at Kozhikode station



It may be noted that the nature of the fourth span (lags 36 to 47) in the PACF and ACF plots is unknown. Therefore, due consideration should also be given for the seasonal autoregressive component  $AR(3)$  (i.e.  $P = 3$ ) and the seasonal moving average component  $MA(3)$  (i.e.  $Q = 3$ ). In the model development phase, it is necessary for the developer to ensure that the model is parsimonious. In order to satisfy the parsimony principle, the sum of the parameters  $p$ ,  $q$ ,  $P$  and  $Q$  of the SARIMA model should be less than or equal to six. Therefore, these four parameters of the candidate model are perturbed in the range of  $-1$  and  $+1$ . It resulted in 15 possible combinations to build the SARIMA model. Out of these, the best model is the one which minimises AICc (corrected Akaike information criteria) and BIC (Bayesian information criteria). The AICc and BIC values for the 15 models are presented in Table 7. The most suitable model that corresponds to the lowest AICc and BIC value is SARIMA(2, 1, 1)(1, 1, 2)<sub>12</sub>. However, it may be noted that the AICc and BIC values of the other four models SARIMA(2, 1, 2)(1, 1, 1)<sub>12</sub>, SARIMA(2, 1, 1)(1, 1, 1)<sub>12</sub>, SARIMA(2, 1, 1)(2, 1, 1)<sub>12</sub> and SARIMA(3, 1, 1)(1, 1, 1)<sub>12</sub> are also closer to the selected model.

Statistical evaluation of the developed models is carried out using the validation dataset. The computed statistical measures are root mean square error (RMSE), mean absolute error (MAE), mean absolute percentage error (MAPE) and mean percentage error (MPE). The results of the statistical evaluation of the models are presented in Table 8. In the statistical evaluation, the SARIMA(2, 1, 1)(2, 1, 2)<sub>12</sub> model did not produce the best results; nevertheless, the results are very close to the models that produced the best results. Therefore, the SARIMA(2, 1, 1)(1, 1, 2)<sub>12</sub> is considered to be

**Table 7** The AICc and BIC values of the SARIMA models for MMAX variable of the Kozhikode station

SARIMA Model	AICc	BIC
SARIMA(3,1,1)(1,1,1) <sub>12</sub>	-2208.86	-2180.66
SARIMA(3,1,0)(2,1,1) <sub>12</sub>	-2176.63	-2148.43
SARIMA(3,1,0)(1,1,2) <sub>12</sub>	-2180.75	-2152.56
SARIMA(3,1,0)(1,1,1) <sub>12</sub>	-2176.52	-2152.33
SARIMA(2,1,1)(2,1,1) <sub>12</sub>	-2208.06	-2179.86
SARIMA(2,1,1)(1,1,2) <sub>12</sub>	-2212.19	-2183.99
SARIMA(2,1,1)(1,1,1) <sub>12</sub>	-2208.52	-2184.33
SARIMA(2,1,0)(2,1,2) <sub>12</sub>	-2173.92	-2145.72
SARIMA(2,1,0)(2,1,1) <sub>12</sub>	-2168.88	-2144.68
SARIMA(2,1,0)(1,1,2) <sub>12</sub>	-2174.36	-2150.16
SARIMA(2,1,0)(1,1,1) <sub>12</sub>	-2169.02	-2148.83
SARIMA(2,1,0)(1,1,3) <sub>12</sub>	-2173.91	-2145.72
SARIMA(2,1,0)(3,1,1) <sub>12</sub>	-2168.35	-2140.16
SARIMA(2,1,2)(1,1,1) <sub>12</sub>	-2210.84	-2182.64

**Table 8** The statistical evaluation results of the SARIMA models developed for the Kozhikode MMAX

SARIMA Model	RMSE	MAE	MPE	MAPE
SARIMA(3,1,1)(1,1,1) <sub>12</sub>	0.846	0.656	-1.36	2.046
SARIMA(3,1,0)(2,1,1) <sub>12</sub>	1.092	0.891	-2.41	2.8
SARIMA(3,1,0)(1,1,2) <sub>12</sub>	1.121	0.92	-2.541	2.9
SARIMA(3,1,0)(1,1,1) <sub>12</sub>	1.117	0.918	-2.551	2.891
SARIMA(2,1,1)(2,1,1) <sub>12</sub>	0.817	0.627	-1.18	1.951
SARIMA(2,1,1)(1,1,2) <sub>12</sub>	0.878	0.675	-1.454	2.11
SARIMA(2,1,1)(1,1,1) <sub>12</sub>	0.868	0.674	-1.451	2.107
SARIMA(2,1,0)(2,1,2) <sub>12</sub>	1.024	0.825	-2.136	2.588
SARIMA(2,1,0)(2,1,1) <sub>12</sub>	1.056	0.856	-2.27	2.685
SARIMA(2,1,0)(1,1,2) <sub>12</sub>	1.089	0.886	-2.402	2.79
SARIMA(2,1,0)(1,1,1) <sub>12</sub>	1.082	0.882	-2.411	2.776
SARIMA(2,1,0)(1,1,3) <sub>12</sub>	1.022	0.823	-2.127	2.581
SARIMA(2,1,0)(3,1,1) <sub>12</sub>	1.13	0.93	-2.582	2.933
SARIMA(2,1,2)(1,1,1) <sub>12</sub>	0.869	0.679	-1.5	2.12
SARIMA(4,1,0)(1,1,1) <sub>12</sub>	1.162	0.962	-2.738	3.037

an appropriate model for forecasting the MMAX variable of the Kozhikode station.

The selected model is also validated using the ACF of the residuals obtained from the fitted SARIMA(2, 1, 1)(1, 1, 2)<sub>12</sub> model to the complete time-series data. The residual plot and the ACF plot are shown in Fig. 5. Ideally, for a model to be absolutely perfect, it is expected to have autocorrelation of residuals close to zero.

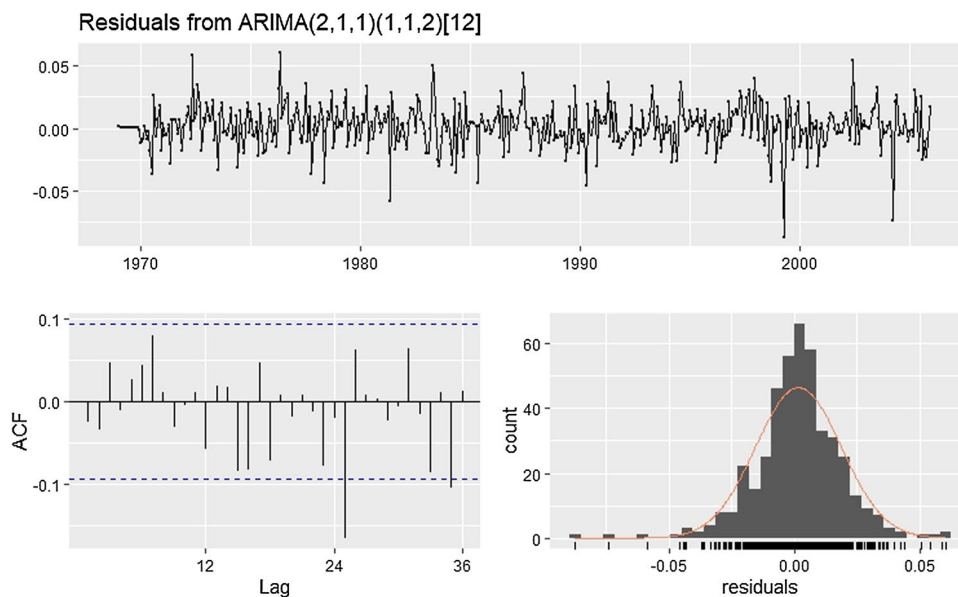
However, if 95% of the spikes lie within the bounds ( $\pm 2\sqrt{T}$ , where  $T$  is the length of the time-series), that would confirm that the series is white noise without autocorrelation. Here, there are two significant spikes (at lag 25 and lag 35), and this translates to 94.4% of the spikes remaining within the bounds. The percentage of spikes lying within the bounds is close to 95% indicates that the selected model has the ability to provide good forecasting results.

Model building and validation for other thirteen variables are carried out using a similar procedure. The data length used for model building and validation is presented in Table 9. Finally, Table 10 shows the apt models for all of the fourteen temperature time-series.

## Summary and conclusions

The development of temperature forecasting models for the state of Kerala, India, is presented in this article. Monthly mean maximum (MMAX) and mean minimum (MMIN) temperature time-series, obtained from seven stations of Kerala is used for the development of the model. The time-series temperature data observed over a period of 47 years, spanning from 1969 to 2015 is utilised in this study.

**Fig. 5** The residual time-series of the fitted SARIMA(2, 1, 1)(1, 1, 2)<sub>12</sub> model for the Kozhikode MMAX variable, ACF plot of the residuals and the distribution of the residuals



**Table 9** The training and the validation data length utilised

Station name	Variable type	Training data	Validation data
Kozhikode	MMAX	1969–2005 (37 years)	2006–2015 (10 years)
Kozhikode	MMIN	1969–2005 (37 years)	2006–2015 (10 years)
Kannur	MMAX	1981–2005 (28 years)	2006–2015 (7 years)
Kannur	MMIN	1981–2005 (28 years)	2006–2015 (7 years)
Alappuzha	MMAX	1969–2005 (37 years)	2006–2015 (10 years)
Alappuzha	MMIN	1969–2005 (37 years)	2006–2015 (10 years)
Punalur	MMAX	1969–2005 (37 years)	2006–2015 (10 years)
Punalur	MMIN	1969–2005 (37 years)	2006–2015 (10 years)
Kottayam	MMAX	1969–2003 (35 years)	2005–2011 (8 years)
Kottayam	MMIN	1969–2003 (35 years)	2005–2011 (8 years)
Thiruvananthapuram	MMAX	1969–2005 (37 years)	2006–2015 (10 years)
Thiruvananthapuram	MMIN	1969–2005 (37 years)	2006–2015 (10 years)
Trivandrum Airport	MMAX	1969–2005 (37 years)	2006–2015 (10 years)
Trivandrum Airport	MMIN	1969–2005 (37 years)	2006–2015 (10 years)

Some data gaps are identified in the datasets obtained from IMD. The missing values are estimated using the expectation–maximisation algorithm. It is natural for a long term time-series dataset of a meteorological variable to possess a trend. Moreover, the monthly mean of the meteorological variable is bound to have seasonal variations. The inherent seasonality in the variable induces a non-stationarity in the time-series datasets. Statistical analysis is carried out on the time-series datasets to understand the nature of the data.

The results from the descriptive statistics indicated that most of the temperature time-series are kurtotic. A preliminary analysis is carried out to test the normality of the data and to check the presence of outliers. The Shapiro-Wilk test and the Grubb’s test are conducted to test the normality and

to check the outliers. The results indicated that the time-series datasets are non-normal and outliers are absent.

The trend analysis is carried out by applying Mann–Kendall’s trend test and Sen’s Slope estimation. The results indicated the presence of trend in at least ten of the fourteen time-series datasets. This served as the first indication for the non-stationary nature of the datasets. In order to confirm the presence of seasonality, with absolute confidence, STL decomposition and KPSS test are conducted. In STL decomposition, the time-series is decomposed into trend, seasonal, and remainder components. The results obtained from these tests clearly indicated the presence of seasonality and thereby, confirmed the non-stationarity of the all the fourteen time-series datasets. Subsequently, one seasonal difference and one first-order difference are applied to

**Table 10** The best-fit SARIMA models developed for forecasting the variables

Station name	Variable type	Number of valid models	Best-fit SARIMA model	AICc	BIC
Kozhikode	MMAX	15	SARIMA(2,1,1)(1,1,2) <sub>12</sub>	−2212.19	−2183.99
Kozhikode	MMIN	31	SARIMA(2,1,1)(1,1,1) <sub>12</sub>	−2026.97	−2002.77
Kannur	MMAX	31	SARIMA(2,1,1)(1,1,1) <sub>12</sub>	−1545.43	−1523.03
Kannur	MMIN	15	SARIMA(3,1,1)(1,1,1) <sub>12</sub>	−1446.43	−1420.34
Alappuzha	MMAX	15	SARIMA(2,1,1)(1,1,1) <sub>12</sub>	−2154.94	−2130.74
Alappuzha	MMIN	31	SARIMA(2,1,1)(1,1,1) <sub>12</sub>	−2045.26	−2021.06
Punalur	MMAX	15	SARIMA(2,1,1)(1,1,1) <sub>12</sub>	−1917.4	−1893.2
Punalur	MMIN	31	SARIMA(2,1,1)(2,1,1) <sub>12</sub>	−1822.05	−1793.86
Kottayam	MMAX	15	SARIMA(2,1,1)(1,1,1) <sub>12</sub>	−1736.01	−1712.17
Kottayam	MMIN	31	SARIMA(3,1,1)(1,1,1) <sub>12</sub>	−1822.03	−1794.25
Thiruvananthapuram	MMAX	31	SARIMA(2,1,2)(1,1,1) <sub>12</sub>	−2154.22	−2126.03
Thiruvananthapuram	MMIN	31	SARIMA(2,1,1)(1,1,1) <sub>12</sub>	−2165.97	−2141.77
TrivandrumAirport	MMAX	15	SARIMA(2,1,2)(1,1,1) <sub>12</sub>	−2341.1	−2312.9
TrivandrumAirport	MMIN	31	SARIMA(2,1,1)(1,1,1) <sub>12</sub>	−2054.48	−2030.29

transform the non-stationary datasets into stationary datasets. The results assist in identifying the values of the differencing parameters necessary for building the SARIMA model.

The SARIMA models are developed individually for each of the fourteen variables using the original time-series datasets. The SARIMA models are developed individually for each of the fourteen time-series datasets. The results indicated that the SARIMA(2, 1, 1)(1, 1, 1)<sub>12</sub> model the ideal one to forecast eight out of the fourteen time-series variables. In order to have a better understanding of local influences, the studies must be carried out on a better spatial and temporal scales.

**Acknowledgements** The first author would like to thank the Ministry of Human Resource Development (MHRD) for the financial support, which was provided in the form of research stipend. The authors also acknowledge the Indian Meteorological Department (IMD) for the temperature datasets.

## Compliance with ethical standards

**Conflict of interest** The authors declare that they have no conflict of interest.

## References

- Aguado-Rodríguez GJ, Quevedo-Nolasco A, Castro-Popoca M, Arteaga-Ramírez R, Vázquez-Peña MA, Zamora-Morales BP (2016) Predicción de variables meteorológicas por medio de modelos arima. *Agrociencia* 50(1):1–13
- Allen MR, Ingram WJ (2002) Constraints on future changes in climate and the hydrologic cycle. *Nature* 419(6903):228–232
- Andronova NG, Schlesinger ME (2000) Causes of global temperature changes during the 19th and 20th centuries. *Geophys Res Lett* 27(14):2137–2140
- Arjun KM (2013) Indian agriculture-status, importance and role in Indian economy. *Int J Agric Food Sci Technol* 4(4):343–346
- Cleveland RB et al (1990) Stl: a seasonal-trend decomposition procedure based on loess. [citeulike-article-id:1435502](https://doi.org/10.1002/citeulike-article-id:1435502)
- Gilbert RO (1987) *Statistical methods for environmental pollution monitoring*. Wiley, New York
- Gocic M, Trajkovic S (2013) Analysis of changes in meteorological variables using Mann–Kendall and Sen’s slope estimator statistical tests in Serbia. *Glob Planet Change* 100:172–182
- Hänsel S, Medeiros DM, Matschullat J, Petta RA, de Mendonça Silva I (2016) Assessing homogeneity and climate variability of temperature and precipitation series in the capitals of North-Eastern Brazil. *Front Earth Sci* 4:29
- Hyndman RJ, Athanasopoulos G (2018) *Forecasting: principles and practice*. OTexts, Melbourne
- Indian Network for Climate Change Assessment and India Ministry of Environment (2010) *Climate Change and India: a 4 × 4 assessment, a sectoral and regional analysis for 2030s, vol 2*. Ministry of Environment & Forests, Government of India, New Delhi
- Jain SK, Kumar V (2012) Trend analysis of rainfall and temperature data for India. *Curr Sci* 102:37–49
- Kang H (2013) The prevention and handling of the missing data. *Korean J Anesthesiol* 64(5):402
- Kendall M (1975) *Rank correlation methods*. Charles Griffin, London (There is no corresponding record for this reference)
- Kocsis T, Kovács-Székely I, Anda A (2017) Comparison of parametric and non-parametric time-series analysis methods on a long-term meteorological data set. *Cent Eur Geol* 60(3):316–332
- Kocsis T, Kovács-Székely I, Anda A (2020) Homogeneity tests and non-parametric analyses of tendencies in precipitation time series in Keszthely, Western Hungary. *Theor Appl Climatol* 139(3–4):849–859
- Kwiatkowski D, Phillips PC, Schmidt P, Shin Y et al (1992) Testing the null hypothesis of stationarity against the alternative of a unit root. *J Econom* 54(1–3):159–178
- Lai Y, Dzombak DA (2020) Use of the autoregressive integrated moving average (ARIMA) model to forecast near-term regional temperature and precipitation. *Weather Forecast* 35:959–976
- Mann HB (1945) Nonparametric tests against trend. *Econom J Econom Soc* 13:245–259
- Mills TC (2014) Time series modelling of temperatures: an example from k efalonia. *Meteorol Appl* 21(3):578–584

- Radziejewski M, Kundzewicz ZW (2004) Detectability of changes in hydrological records/possibilité de détecter les changements dans les chroniques hydrologiques. *Hydrol Sci J* 49(1):39–51
- Sen PK (1968) Estimates of the regression coefficient based on Kendall's tau. *J Am Stat Assoc* 63(324):1379–1389
- Tiwari P, Kar S, Mohanty U, Dey S, Kumari S, Sinha P (2016) Seasonal prediction skill of winter temperature over North India. *Theor Appl Climatol* 124(1–2):15–29
- Trenberth KE (1999) Conceptual framework for changes of extremes of the hydrological cycle with climate change. In: *Weather and climate extremes*. Springer, Dordrecht, pp 327–339
- Wang H, Huang J, Zhou H, Zhao L, Yuan Y (2019) An integrated variational mode decomposition and arima model to forecast air temperature. *Sustainability* 11(15):4018
- Wanishsakpong W, Owusu BE (2020) Optimal time series model for forecasting monthly temperature in the southwestern region of Thailand. *Model Earth Syst Environ* 6(1):525–532





# RS- and GIS-based modeling for optimum site selection in rain water harvesting system: an SCS-CN approach

Khalid Mahmood<sup>1</sup> · Ansab Qaiser<sup>2</sup> · Sumar Farooq<sup>2</sup> · Mehr un Nisa<sup>2</sup>

Received: 10 April 2020 / Accepted: 27 June 2020 / Published online: 10 July 2020  
© Institute of Geophysics, Polish Academy of Sciences & Polish Academy of Sciences 2020

## Abstract

In this study, an integrated approach has been adopted for optimum selection of locations for rain water harvesting (RWH) in Kohat district of Pakistan. Various thematic layers including runoff depth, land cover/land use, slope and drainage density have been incorporated as input to the analysis. Other biophysical criteria such as geological setup, soil texture and drainage streams characteristics were also taken into account. Drainage density and slope were derived from digital elevation model, and map of land use/land cover was prepared using supervised classification of multi-spectral Sentinel-2 images of the area. Aforementioned thematic layers are assigned respective weights of their importance and combined in GIS environment to form a RWH potential map of the region. The generated suitability map is classified into three potential zones: high, moderate and low suitability zones consisting of area 638 km<sup>2</sup> (21%), 1859 km<sup>2</sup> (62%) and 519 km<sup>2</sup> (17%), respectively. The suitability map has been used to mark accumulation points on the down streams as potential spots of water storage. In addition, site suitability of artificial structures for RWH consisting of farm ponds, check dams and percolation tanks has also been assessed, showing 3.2%, 3% and 4.5% of the total area as a fit for each of the structure, respectively. The derived suitability will aid policy makers to easily determine potential sites for RWH structures to store water and tackle acute paucity of water in the area.

**Keywords** Site suitability · Rain water harvesting · Remote sensing · Geographical information system

## Introduction

Water is indeed one of the primary driving forces of our very nature. It is the basic need for human as well as for the animal and plant life. The global requirement of water is intensifying with time due to rapid increase in the world population, modern and improving living standards, industrialization and irrigated agriculture (Buraihi and Shariff 2015). Especially under the current climate change scenario where changed rainfall patterns are causing either scarcity or floods, the natural water balance has been lost. So at this time it is more important to manage fresh water resources than ever it was needed before. South Asian countries are more likely to get affected by this scarcity. Pakistan is a

typical example of a country facing effects of this natural imbalance. It was a water surplus country once but now turning into a water deficit region as the fresh water is depleting rapidly. Pakistan is ranked third among the countries which are facing substantial water crisis unheeded by the authorities (Nabi et al. 2019). And the country may run dry by 2025 if such trend continues. It is an agricultural country and its economy relies heavily on the growth of agriculture sector which needs the water most (Ahmed et al. 2007).

Over the time people have acquired many alternative ways to get water for different purposes including digging wells, harvesting rain water, melting snow and ice, accumulating fog and dew, collecting water from the evapotranspiration of plants (Kadam et al. 2012; Tumbo et al. 2012). One of these management solutions is the rain water harvesting (RWH). In areas where water supply is barely sufficient this solution proves to be a propitious method to support scarce water resources of the area to satisfy the demand (Buraihi and Shariff 2015; Mugo and Odera 2019). This research aims to locate potential areas to efficiently harvest rainwater. At one side, proper management of rain water can help in avoiding

✉ Khalid Mahmood  
khalid.spsc@pu.edu.pk

<sup>1</sup> RS and GIS Group, Department of Space Science, University of the Punjab, Lahore, Pakistan

<sup>2</sup> Department of Space Science, University of the Punjab, Lahore, Pakistan

or at least reducing intensity of floods and at the other end provide naturally filtered water in the dry spells to fulfill domestic and other uses (Sekar and Randhir 2007; Buraihi and Shariff 2015). RWH is a combination of versatile and resourceful techniques to filter, stock, and distribute rain water for various domestic purposes. This is the simplest and easily accessible alternative of water management in areas with sufficient rain as it yields extra water to deal water deficit problems (Helmreich and Horn 2009; Gavit et al. 2018). This process of RWH simply expresses the undeviating hold up of rainwater as surface runoff. For domestic purpose, the surface runoff from the roof of an individual house or from paved surface can be harvested. There are many advantages of RWH, i.e., rainwater is not chlorinated so it is unpolluted and open source of water, this harvested water is ideal for crop planting due to its clarity (Buraihi and Shariff 2015). The degree of harvesting can be enhanced from an individual to a bigger catchment or reservoirs for public use. This process of RWH is also useful for the enhancement of ground water recharge. This system has been proved to be very beneficial for many countries in the world.

Although RWH has its history expanded over centuries but now in these modern days, for better results, before installation of RWH system, there are proposed analyses to check the land suitability for such an installation (Mugo and Odera 2019). The more effective and popular analyses made use of the geographic information system (GIS) and satellite remote sensing (SRS) utilities (Buraihi and Shariff 2015). To serve the purpose, various procedures are in practice. One of the methods is the Analytic Hierarchy Process (AHP), a multi-criteria decision analysis. It assigns weights to each of the input criteria showing its level of contribution in the decision support system. Another GIS-based method used in order to delineate these sites for RWH is the weighted overlay of geographic distribution of the involved input parameters, i.e., drainage density, slope, runoff depth, the soil map and the land-use/land-cover (LULC) map (Kadam et al. 2012; Buraihi and Shariff 2015). It also assigns weight to each of the parametric layers. This study has made use of these SRS- and GIS-based strategies to proposed potential rainwater harvesting sites for the district Kohat, Pakistan.

SRS has been emerged as an alternative of traditional in situ sampling methods that were expensive, time-consuming and tedious in their (Mahmood et al. 2017a; Manzo et al. 2017). Its ability to provide bird eye view of larger area with detailed topography and many other proxy factors helping in understanding ongoing process and natural settings of a region at once are the factors making it a better substitute in many environmental related studies (Manzo et al. 2017; Mahmood et al. 2019). The basic data of SRS is reflectance of Earth's surface measured in various spectral range which is interpreted using spatial analysis of various types provided by GIS (Yan et al. 2014; Manzo et al. 2017). In addition to

processing of SRS data GIS also provide ease and accuracy of many other space related data handling, a typical example of it is the Weighted Linear Combination (WLC) of multiple geographic datasets. A general recommendation by researchers is the SRS data with better spatial resolutions, i.e., QUICKBIRD with pixel dimensions of 0.65 m. However, depending upon phenomenon under consideration, freely available SRS data of Sentinel 2 with generalization dimensions of 10 m may prove to be a very suitable option. So both these techniques (SRS and GIS) are helpful for studying surface phenomenon, i.e., RWH of an area.

### Study area

The study area for this research is Kohat district located in the province Khyber Pakhtunkhwa (KPK) of Pakistan, situated adjacent to the Potwar plateau. The zone lies in the range from 33.06° N to 33.75° N and from 71.06° E to 72.01° E. Geographical association of the study area is shown in Fig. 1. With an area of about 2987 km<sup>2</sup>, it contains a population of 723,000 (Population and Household Detail from Block to District Level: Khyber Pakhtunkhwa, 2018). Administration wise the area has been divided into two tehsils Kohat and Lachi.

The climate of the area is semiarid and sub-humid subtropical continental highland. The average temperature of the region varies with altitudes, i.e., plains are relatively warmer and mountains are cooler. The weather remains hot from May to September with peak in June (average maximum temperature 41 °C), whereas coldest month is January (average minimum temperature 5 °C). Annual average temperature in the region is around 24 °C. Monthly average temperature of the region, as per 40 years' record from 1978 to 2017, is shown in Fig. 2. Annual average precipitation in the region, as per 40 years' record from 1978 to 2017, is 580 mm. The monsoon downpours peak in July and August. Looking at the seasonal rain fall patterns (Fig. 3), Kohat district has an additional advantage of two well separated peaks of higher precipitation. One of the peak, relatively smaller, is in March with 82 mm precipitation in a month, for which the surplus rainwater can serve as a reserve for the dry months of May and June, whereas the major peak centers in July and August with about 180 mm in 2 months, for which the collected rainwater can serve as a source for the upcoming dry months of October, November and December. So this way Kohat has two sets of supply and demand periods and seasonally the RWH system of the region can be divided into two temporal frames leading toward smaller collection units with maximum efficiency.

Fig. 1 Study area map

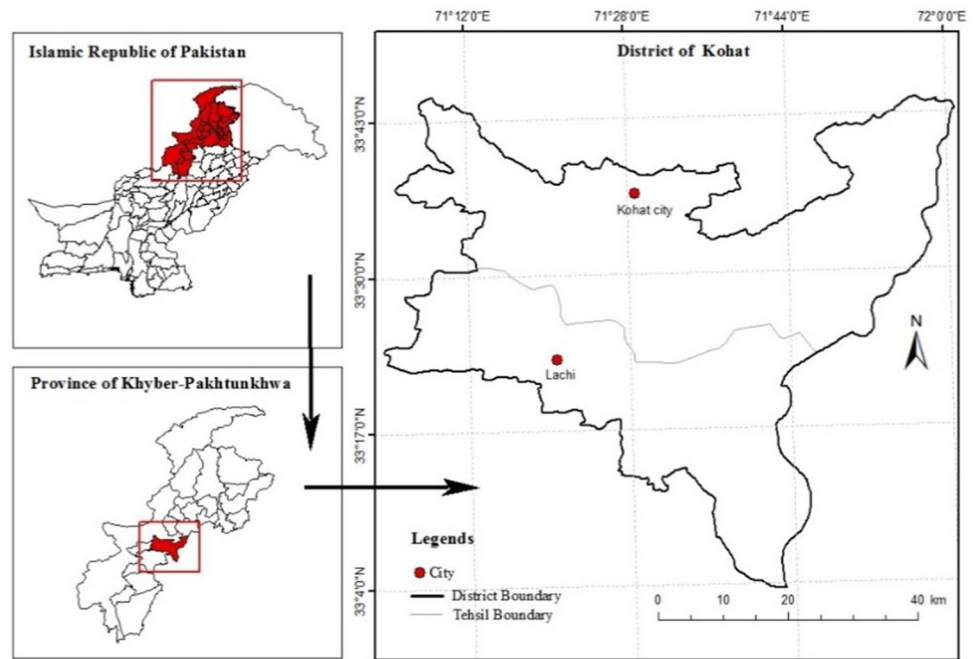
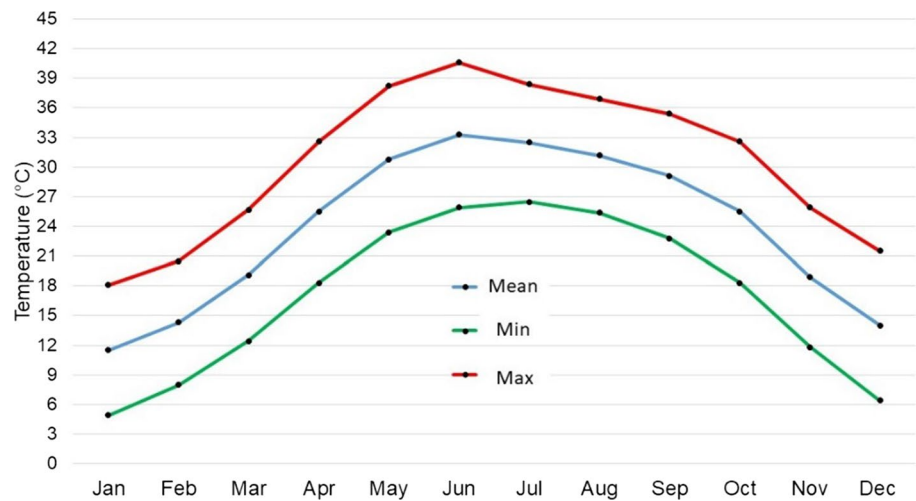


Fig. 2 Average monthly temperature profile of Kohat (from 1978 to 2017)



## Materials and methods

### Datasets collection

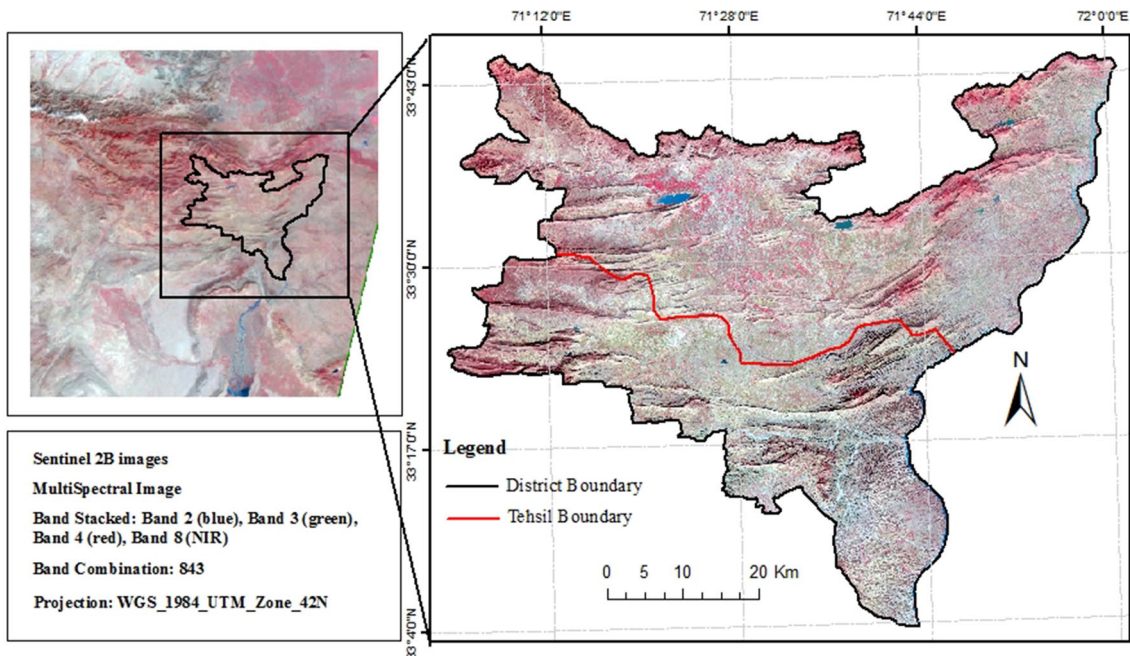
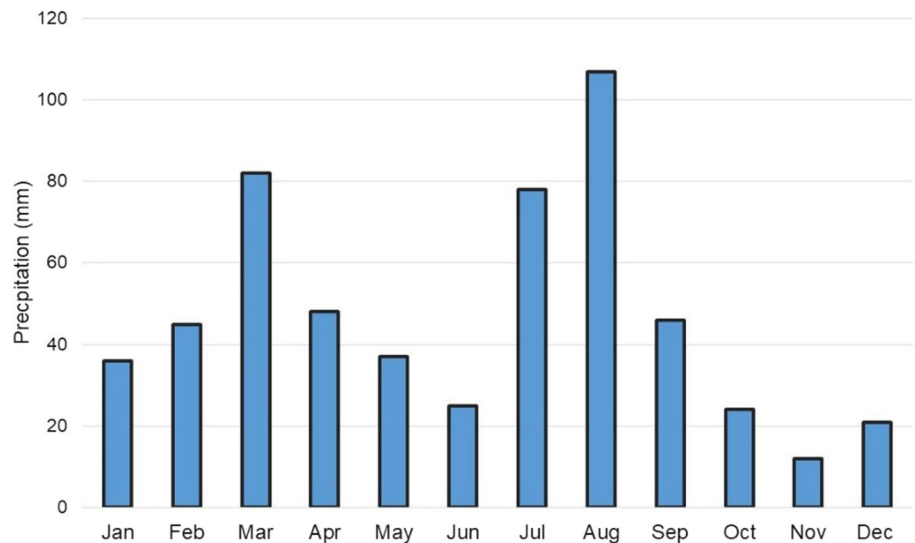
The SRS data of Sentinel 2B acquired on 15-11-2018 has been used for preparing LULC of the area. For elevation, data of SRTM DEM with spatial resolution of 30 m has been used. A detailed soil map, topographic information, soil texture triangle and soil series information of Kohat district were acquired from Soil Survey of Punjab. Finally, for accurate rainfall information, data of past 10 years from two sources were obtained, Pakistan Meteorological

Department (PMD) and Soil Survey of Punjab (Soil Survey of the Punjab, Pakistan).

### Datasets preparation

The preprocessing on acquired four Sentinel-2B images was done, during the preprocessing stage, the bands (red, green, blue, NIR) with 10 m spatial resolution were stacked together, the four images were mosaicked, and the area of interest was clipped using ERDAS Imagine platform; the output is shown in Fig. 4. To prepare a Hydrological Soil Group (HSG) map the geo-referencing and digitization of

**Fig. 3** Average monthly rainfall of Kohat (from 1978 to 2017)

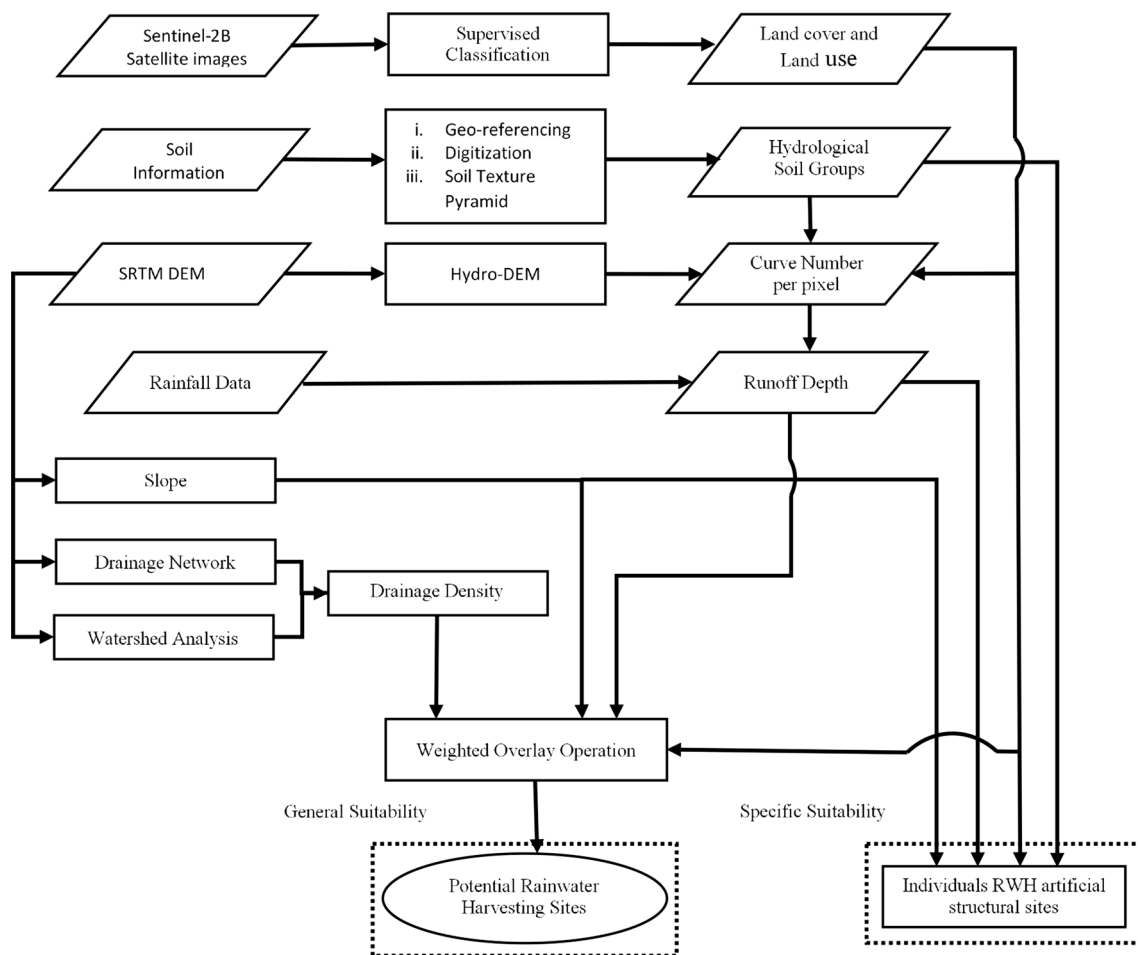


**Fig. 4** Sentinel 2B image (Kohat district)

detailed soil information map of Kohat was carried out using Arc GIS as the GIS platform. Later on, hydrological soil groups were computed using soil texture triangle and Kohat soil series information. The 40 years' rainfall information of Kohat district was utilized to calculate average monthly as well as annual rainfall of the area, the monthly average details are given in Fig. 3, and annual average is 560 mm.

Overall layout of the methodology has been given in Fig. 5. Solution of the problem statement initiated with the preparation of thematic layers using the collected data sets. This selection of the four contributory thematic layers

comprises of runoff depth, slope, drainage density and land use/land cover which has been made based on the reviewed literature (Khalid et al. 2017). In order to prepare a land-use/land-cover layer of the region, preprocessing and classification of Sentinel-2B images performed in ERDAS Imagine, followed by supervised classification using maximum likelihood algorithm of images to six classes, were made, i.e., bare land, grassland, crops, fallow, urban, and open water. Topographic slope layer has been prepared by selecting appropriate z-factor value which computed to be 0.00001036 depending on latitude range of the area by method devise by



**Fig. 5** Flow chart diagram of methodology

Fry (2007). Drainage density is defined as “the total length of streams per unit catchment area” and mathematically as follows (Dragicevic et al. 2019).

$$DD = \frac{\sum_{i=1}^n L}{A} \quad (1)$$

where  $n$  is number of streams,  $L$  is length of the streams (km) and  $A$  is contributing drainage area (km). For quantifying drainage density, drainage network information and watershed delineation were required. Drainage network was derived from DEM using stream raster and flow direction raster, and then “Stream to Feature” tool was used to convert it to vector format. Watershed delineation was carried out using flow-direction raster and stream-link raster, using hydrology tools. Spatial join has been used to combine the information of area of watershed and lengths of the streams falling in the respective watersheds for calculating drainage density per watershed.

Assessment of surface runoff has been made using SCS-CN method, initially developed by USDA, Soil

Conservation Service (SCS). It is explained in fourth section of National Engineering handbook (NEH) (Ponce and Hawkins 1996). The main cause of success of SCS-CN is that this method compiles the parameters that affect the generation of runoff like, soil types, land use and land cover, moisture conditions of that area, surface condition, incorporated by single CN variable (Souliis et al. 2009; Kadam et al. 2012). The method based on the calculation of water balance is written as follows (Li et al. 2015):

$$P = I_a + Q + F \quad (2)$$

$$\frac{Q}{P - I_a} = \frac{F}{S} \quad (3)$$

where  $P$  is total rainfall (mm),  $I_a$  is initial abstraction (mm),  $F$  is cumulative infiltration excluding  $I_a$  (mm),  $S$  is potential maximum retention (mm) and  $Q$  is direct runoff (mm). By the combination of Eqs. (2) and (3), standard form of SCS-CN method turns out into following equation.

$$Q = \frac{(P - Ia)^2}{P - Ia + S} \tag{4}$$

which is effective when  $P \geq Ia$ ; else  $Q = 0$ . This approach relies on two essential assumptions; first ratio of maximum possible runoff to actual rainfall is equal to ratio of real infiltration to maximum possible retention (Satheeshkumar et al. 2017). As per the second supposition, the volume of  $Ia$  is the segment of maximum possible retention.  $Ia = 0.2S$  (Li et al. 2015; Satheeshkumar et al. 2017).

$$Q = \frac{(P - 0.2S)^2}{P + 0.8S} \tag{5}$$

$S$  is calculated by using a mathematical mapping equation depicted in the form CN as follows:

$$S = \frac{25400}{CN} - 254 \tag{6}$$

where CN is the curve number that depends on LULC, HSG and AMC (Antecedent Moisture Condition) and can be obtained from SCS handbook of Hydrology (NEH-4), section-4 (Satheeshkumar et al. 2017). In addition to that Arc Hydro extended tool has a built-in function to generate CN lookup table and CN value raster. It has no dimensions and has a range of 0-100 and depicts the abstraction properties of watershed. Ideal Impermeable surfaces such as water surfaces where all rainfall become runoff would have  $CN = 100$ . And the surfaces which absorb all rainfall would have  $CN = 0$  (Gray and Burke 1983).

Runoff can be computed using Eq. (5), provided the value of CN is known. Estimating the CN for a catchment is considered an important application of GIS. In this research HEC-GeoHMS (Geospatial Hydrologic Modeling Extension) incorporated by Arc Hydro Tool has been used to investigate the value of curve number raster. Hydrological soil group chart has been prepared and is shown in Table 1. The participating layers have been classified before their unification using weighted overlay analysis (Satheeshkumar et al. 2017).

Hydrological Digital Elevation Model (Hydro-DEM) was generated from Arc hydro extension of ArcGIS. CN-Look-up table was made using curve number details with respect

**Table 2** Factors and scale values of different factors (Buraihi and Shariff 2015)

Factor	Weight of class (Pi) (%)	Classes	Rank of class (Wi)
Land use	7	Barren	9
		Grassland	7
		Crops	3
		Fallow	3
		Urban	1
		Water	1
		Slope (%)	30
	4.28–10.36	9	
	10.37–17.98	7	
	17.99–27.73	5	
	27.74–77.7	1	
Runoff (mm)	48	217–301	1
		302–411	3
		412–462	5
		463–479	7
		480–529	9
		Drainage density (km/km <sup>2</sup> )	15
	0.15–0.37	3	
	0.38–0.57	5	
	0.58–0.86	7	
	0.87–1.87	9	

to land cover as defined by USDA. Along with LULC and HSG merged layer, hydro-DEM and CN-Look-up table, “Generate CN Grid” tool from HEC-GeoHMS was used in the estimation of curve numbers per pixel (Amakrishnan et al. 2009; Shukur 2017).

The weights, showing relative importance of each of the parameters in assessing RWH potential, need to be specified so that contributing factor of each of used parameters can be controlled. For this study these weights to each of the contributing factors have been assigned based on the reviewed literature, followed by pairwise comparison metrics analysis (Maina and Raude 2016; Mugo and Odera 2019). Classification of all the input variables along with their weight of importance has been given in Table 2. Finally, weighted sum

**Table 1** Soil Conservation Service classification

Hydrologic soil (HSG)	Soil textures	Runoff potential	Final infiltration
Group A	Deep, well-drained sands and gravels	Low	> 7.5
Group B	Moderately deep, well drained with Moderate	Moderate	3.8–7.5
Group C	Clay loams, shallow sandy loam, soils with moderate to fine textures	Moderate	1.3–3.8
Group D	Clay soils that swell significantly when wet	High	< 1.3

has been carried out to unified score for each of the location using individual ranks assigned by each of the parameters.

### The specific suitability of individual RWH storage structures

In order to find suitability sites for individual RWH storage structures, the generated thematic layers of LU/LC, slope, HSG and runoff were considered. The layers were overlaid and suitable sites for RWH structures like farm ponds, check dams and percolation tanks were found. Characteristics of each criterion on the basis of which the sites for each structure were found are given in Table 3.

## Results and discussion

This study has made use of different parameters. and each of them has its critical role in deciding ability of a location to be a potential RWH site. Out of all the four layers, runoff depth was assigned the highest weight (48%), whereas the lowest weight was assigned to LULC that is (7%). The slope and drainage density were given 30% and 15%, respectively. Assessing land cover distribution in the region, 0.65% out of the total land cover is water majorly consisting of small lakes (spatially found concentrated in the north) and stream/ rivers flowing along eastern boundary of the area. This small percentage is considered to be absolutely unsuitable for the installation of any rainwater harvesting system. Grasslands occurring at the outskirts along northern and western boundary has a coverage of 29% and 0.5% is covered by crops and lies in the northeast, whereas the urban settlements are occupying a small percentage of 0.2% and definitely unsuitable for constructing RWH sites. Fallow lands are peppered throughout the study area, covering 11%, whereas major portion (59%) is barren land. Spatial arrangement of all these land covers is shown in Fig. 6a.

The soils of Kohat region were categorized into four hydrological sets A, B, C and D in accordance with the rates of infiltration of different types of soils on the basis of their textures, for example loamy, sandy clay loam, etc., as referred by Soil Conservation Service Classification (USDA 1974). The type A soil (the well-drained soils and gravels), mainly covering northern regions of the area, consists

of about 135 km<sup>2</sup>. Type B soil (the moderately drained), expanding toward north from center, has an area coverage of about 438 km<sup>2</sup>. Type C soil (clay loams), having coverage along eastern and western boundaries of the area, is found to have an area value approximately 364 km<sup>2</sup>. The Type D is with the maximum area coverage of about 2050 km<sup>2</sup> and is situated mostly on edges of the region. As it is previously mentioned in Table 1, the suitable soil type for installing RWH structure is type D because of its high runoff potential and low infiltration rate and luckily it is the top existing soil (70% of study area). The spatial distribution of these four hydrological soil groups is shown in Fig. 6b.

An average rainfall of 529 mm is generating surface runoff ranging from 217 mm to 529 mm depending on the geographical situations, whereas for a site to be suitable the runoff should be greater than 300 mm so 98% of the area satisfying this condition (Buraihi and Shariff 2015), while the remaining 2% area of non-suitability is covered either by urban settlements or water bodies, which is already excluded for the potential list. Figure 6c illustrates spatial distribution of curve numbers on the basis of which runoff depth has been calculated that is shown in Fig. 6d. Slope was assigned the second highest weight (30%) in the analysis and ranges between 0%–77%, as shown in Fig. 6e. Most of the areas have range of 0–10% that is considered RWH site, while the edges of the boundary in the north part of Kohat are considered to be unsuitable for rainwater harvesting due to the very high values of slope. Drainage density, carrying an importance of 15% in the analysis, has been shown in Fig. 6f for individual watersheds that ranges from 0 to 1.874 per km which is quite suitable for potential RWH site.

### RWH potential suitability map

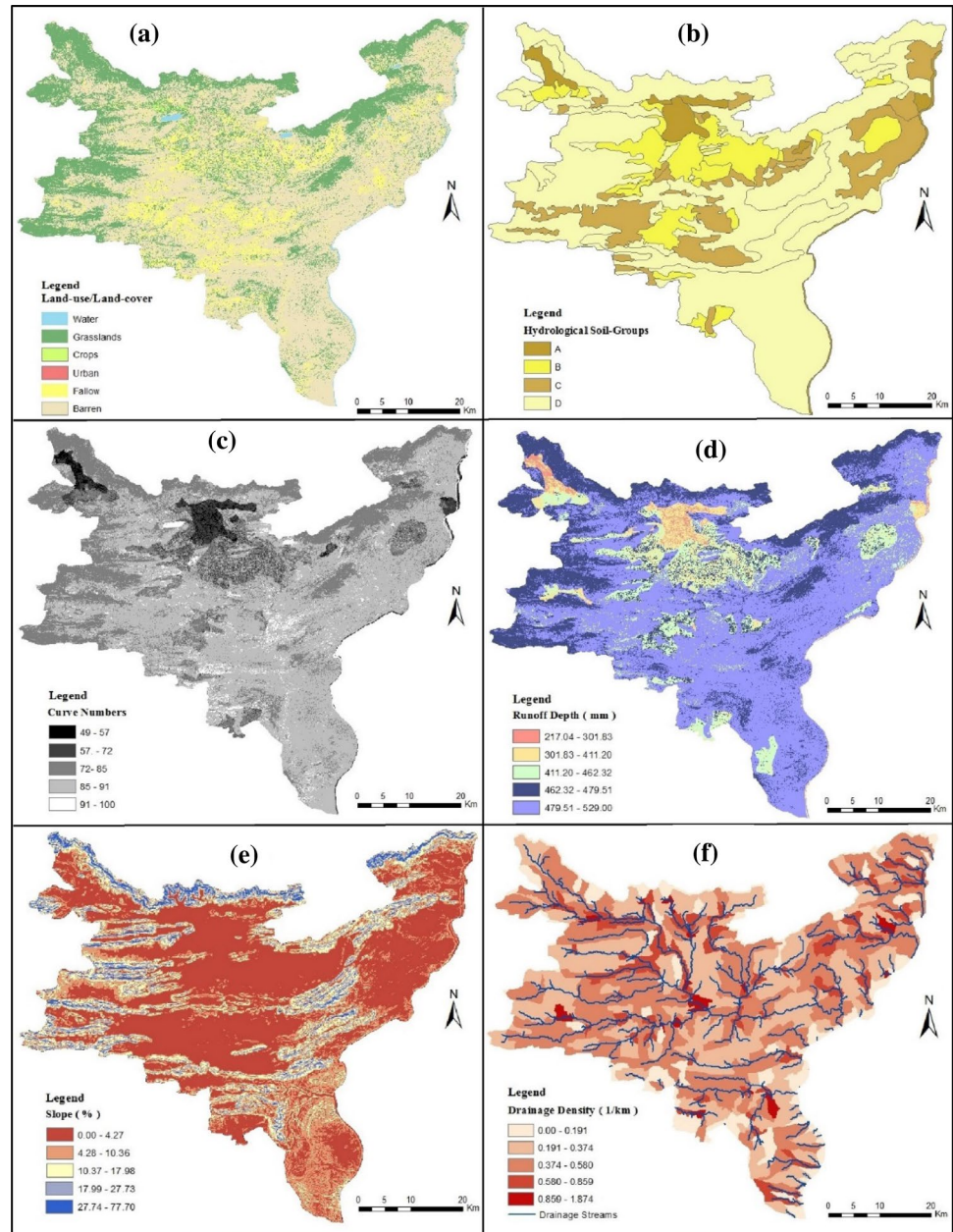
The final output of unified suitability score with geographical spread has been given in Fig. 7. The output has been divided into three classes of suitability level and percentage area for each of the class has been calculated to assess overall suitability of the study area for RWH potential. This classification over weighted overlay raster has been performed using equal interval classifiers that divides input values to specified number of classes (three in this particular case) while giving same value intervals to each of the class.

Only 17.2% of the study area has been categorized into poorly suitable as potential site for RWH. This area mainly consists of urban settlements and the water bodies which is already discussed as the non-suitable land covers for such an arrangement. Major portion (61.6%) of the study area lies with optimized score range of moderately suitability and 21.1% has been assessed as the region of top suitability. For the most suitable sites, the runoff values were maximum, the infiltration rate was least, and land cover was mostly the barren land. However, these regions of high suitability have

**Table 3** Selection of artificial structures for RWH (Ammar et al. 2016; Khalid et al. 2017)

Artificial RWH structure	Slope (%)	Land cover type	Soil type
Farm-pond	< 5	Agriculture	C and D
Check dam	< 15	Barren	C
Percolation tank	< 10	Barren	B

**Fig. 6** Geographical display of various parameters (**a** LULC; **b** hydrological soil groups; **c** curve numbers; **d** runoff depth; **e**: slope; **f**: drainage density)



their spread throughout the region, with high concentration in the extreme southern parts which can get the maximum accumulation of rain water. Similarly, areas with low suitability have high concentration in the north. Emergence of these patterns with large rainwater collecting areas as the most suitable sites is an advantage of the region, making the area naturally blessed with high potential of success rate as a pilot project of RWH. Looking at the potential points of rain water collection, the highly suitable stream points are distributed throughout the area, not showing the similar extreme concentration in the south as was shown by the suitable geography. It is because of the fact that the points have been marked at end of the suitable stream with some specific

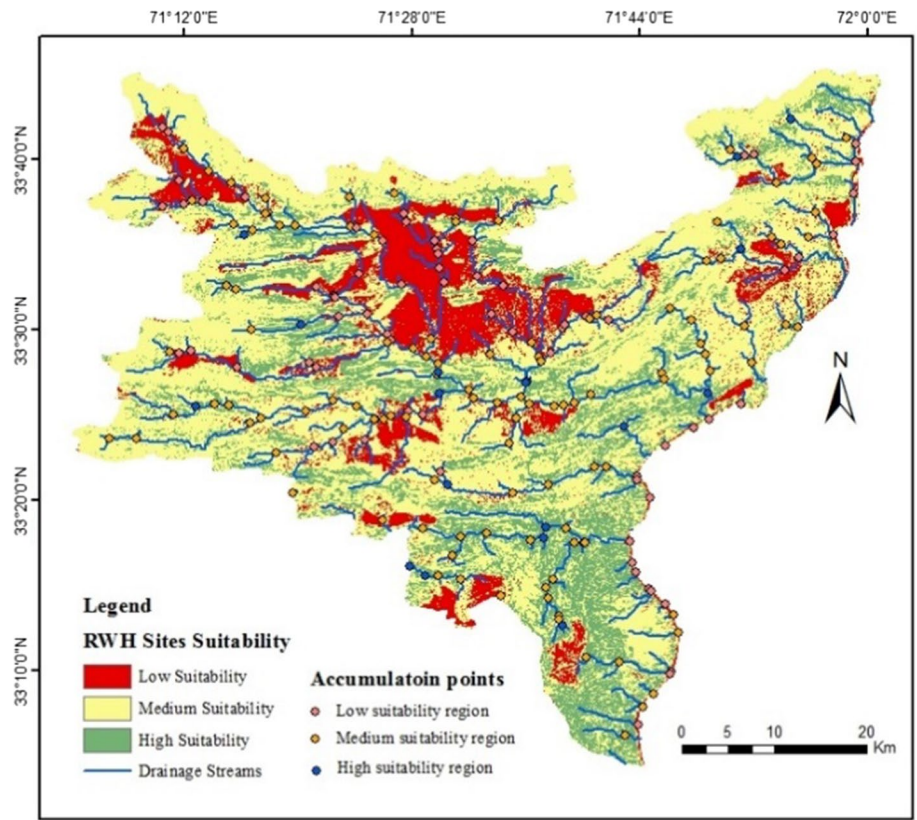
amount of water to be tackled at the location, whereas the low suitability points are well concentrated in the low suitable geography in the central north. The distribution of most suitable location throughout the area can lead to a highly efficient system of RWH in the region in terms of water storage as well as utilization management.

### Artificial storage structures map

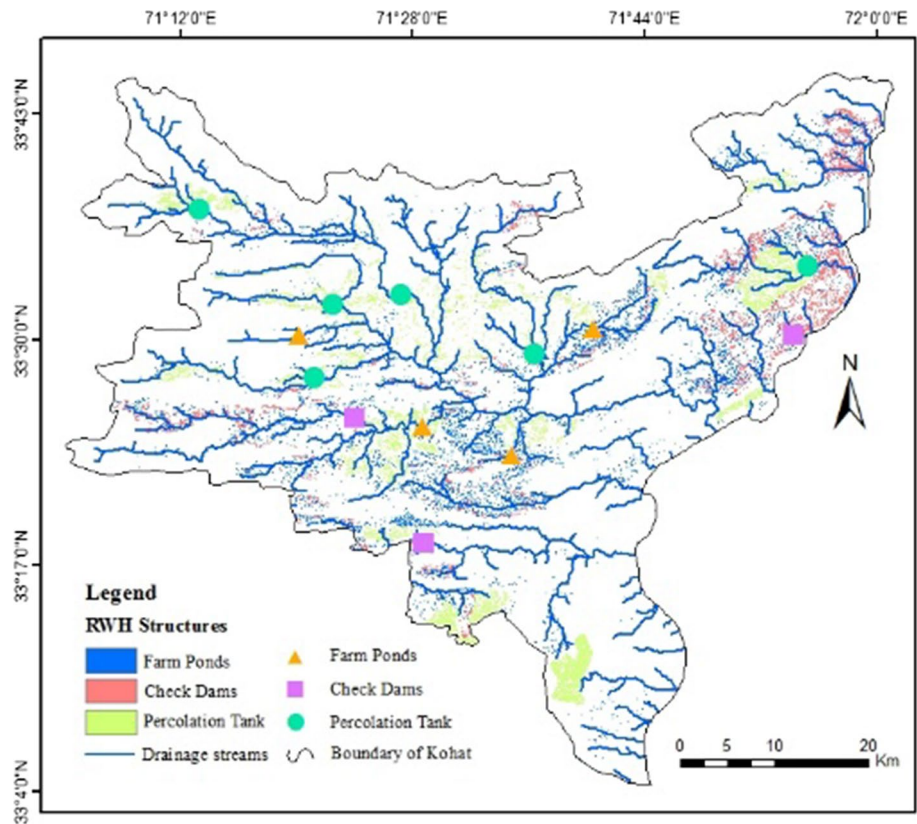
In addition to general suitability of sites, additional analysis has been performed to assess sites for artificial RWH storage structures, i.e., check dams, farm ponds and percolation tanks as shown in Fig. 8. It will prove a further help



**Fig. 7** Spatial distribution of RWH potential sites and suitable accumulation points



**Fig. 8** Artificial structures for RWH



for development of application specific infrastructure to the decision makers in the region. These three types of structures need specific soil type and runoff depth.

Though the critical analysis 3.21% of the study area is found suitable for farm ponds that needs slope range of 0–5%, runoff depth lies from moderate to high and soil type as sandy clay loam. About 3% of the study area satisfies the conditions for the development of check dams with slope between 5 and 15% and runoff depth and soil type same as for farm ponds. The third structure is percolation tank with a suitability area of about 4.5% of the total, which can be installed in the regions where slope value is less than 10%, runoff is low and soil type is clay. Over this area allocation, the possible locations of individual structures, with maximum benefits, have also been pointed out and are shown in Fig. 7.

Well-distributed six locations have been marked for percolation tanks, four for the farm ponds concentrated mainly in the center of the region, and three for check dams. Installation of percolation tanks has been marked with four points at upstream of the main residential region of Kohat with two points as a facility to relatively low populated areas in the east and west of it. All the farm ponds suggestions have been emerged from periphery of the residential area where most of the farm land and farming business is concentrated. Three check dam suggestions are also providing a perfect distribution in the downslope of main Kohat and will be of great benefit for southern settlements to this main population zone. Overall the arrangement of potential sites for RWH system has come out with maximum potential benefits providing the use techniques as the perfect solution for planning such a system and the chosen area as the perfect fit of a pilot project for the RWH.

## Conclusion

The aim of this study was to identify the suitable sites for harvesting rain water in the region of Kohat, Pakistan. The study area has been found having a natural setup to support RWH as an alternate for fulfilling water requirements at different scales. If the rainwater is harvested, filtered, stored and managed properly, the water scarcity problem can be fixed to considerable extent in the area. This research successfully develops the suitable sites in the area where rainwater can be harvested with the installation of different artificial RWH structures. Geographical multicriteria evaluation process has been successfully utilized for analyzing site suitability that is based on various methods of GIS and SRS. These methods have been put together in such a sequence that suits with natural settings of the phenomenon. The used thematic layers in the analysis are slope, land cover/land use, runoff depth, and drainage density that were combined using

weighted overlay process in GIS environment. The study has provided many insights into the spatial distribution of RWH controlling parameters.

About 0.65% of the study area is covered by water that is absolutely not suitable for building any RWH structure and major land cover is barren that comprises about 58% of the area. The most appropriate soil type has a spread, covering of 70% and 98% of the area having a runoff greater than 300 mm, needed for building a RWH system. Similarly, maximum of the area has a suitable slope value of 0–10% and drainage density of all the individual watersheds is also fit for planning a RWH system. Although individual parameters have suitable scenario over major portion of the area, a geographical mismatch of their combination has somehow decreased this area to suitability. Even then the final suitability map has come out with only 17% of the area that is not suitable. Still 83% of the study area is supporting development of RWH system in the region with 21% area as highly suitable. So overall the study area is naturally blessed with all the ingredients needed for the development of sustainable and efficient RWH systems.

The area in immediate surroundings of Kohat city lies in the low suitable region while most of the highly suitable sites lie in the center and the southern parts. Conclusively, the generation of accumulation points at the end of down streams proved that there are enough points of accumulation throughout the region where different artificial recharge/storage structures can be installed in the selected sites which can prove to be greatly helpful in conserving water in water deficit areas.

In addition to this general suitability of RWH, sites for specific structure have also been determined. The area of suitability found for farm ponds is 3.2% with four construction sites, for check dams it is 2.9% with three construction sites and for percolation tanks it is 4.5% of the total with six construction sites. Overall the arrangement of potential sites, both general and specific, for RWH system in the area is proving beneficial in two ways. At first, emergence of these patterns with large rainwater collecting areas as the most suitable sites and distribution of these resource points throughout the area both are leading to highly efficient system in terms of water storage and utilization. This is an advantage of the region, making the area naturally blessed with high potential of success rate as a pilot project of RWH. At second place perfection in the output with maximum potential benefits has proved the used techniques as the preferred solution for planning such a system.

The SCS-CN approach adopted for the study has been proved to be a better technique not only for determining best sites for RWH system but also to assess natural settings of an area for the purpose. This approach can be implemented as it is to almost all the hilly areas. As the approach provides pre-assessment of RWH potential so when combined

with some existing ranking criteria like one proposed by Mahmood et al. (2017b) can lead toward a better and beneficial resources allocation for such systems in the world.

**Authors contribution** All the authors are genuine contributor of the research work.

**Funding** There is no funding for this research.

**Availability of data and material** The used data is available with Soil Survey of Pakistan, Lahore, Pakistan, meteorological data with Pakistan Meteorological Department and used satellite images with Copernicus Open Access Hub. Additionally, all this data is also available with authors.

## Compliance with ethical standards

**Code availability** Not applicable.

**Conflict of interest** There is no conflict of interest.

## References

- Ahmed A, Iftikhar H, Chaudhry GM (2007) Water resources and conservation strategy of Pakistan. *Pak Dev Rev* 46(4):997–1009
- Amakrishnan DR, Andyopadhyay AB, Kusuma KN (2009) SCS-CN and GIS-based approach for identifying potential water harvesting sites in the Kali Watershed, Mahi River Basin, India. *J Earth Syst Sci* 118:355–368
- Ammar A, Riksen M, Ouassar M, Ritsema C (2016) Identification of suitable sites for rainwater harvesting structures in arid and semi-arid regions: a review. *Int Soil Water Conserv Res* 4:108–120
- Buraihi FH, Shariff ARM (2015) Selection of rainwater harvesting sites by using remote sensing and GIS techniques: a case study of Kirkuk, Iraq. *J Teknol* 76(15):75–81
- Dragicevic N, Karleusa B, Ozanic N (2019) Different approaches to estimation of drainage density and their effect on the erosion potential method. *Water* 11:593
- Fry C (2007) Setting the Z factor parameter correctly, imagery & remote sensing. ESRI. <https://www.esri.com/arcgis-blog/products/product/imagery/setting-the-z-factor-parameter-correctly/>. Accessed 12 June 2007
- Gavit BK, Purohit RC, Singh PK, Kothari M, Jain HK (2018) Rainwater harvesting structure site suitability using remote sensing and GIS. *Hydrologic modeling*. Springer, Singapore, pp 331–341
- Gray DD, Burke CB (1983) Occurrence probabilities of antecedent moisture condition classes in Indiana. Report, Purdue University, Indiana
- Helmreich B, Horn H (2009) Opportunities in rainwater harvesting. *Desalination* 248(1–3):118–124
- Kadam AK, Kale SS, Pande NN, Pawar NJ, Sankhua RN (2012) Identifying potential rainwater harvesting sites of a semi-arid, basaltic region of Western India, using SCS-CN method. *Water Resour Manag* 26(9):2537–2554
- Khalid J, Marsumi A, Shamma AMA (2017) Selection of suitable sites for water harvesting structures in a flood prone area using remote sensing and GIS-case study. *J Environ Earth Sci* 7(4):91–100
- Li J, Liu C, Wang Z, Liang K (2015) Two universal runoff yield models: SCS vs. LCM. *J Geogr Sci* 25:311–318
- Mahmood K, Batool A, Faizi F, Chaudhry MN, Ul-Haq Z, Rana AD, Tariq S (2017a) Bio-thermal effects of open dumps on surroundings detected by remote sensing-influence of geographical conditions. *Ecol Ind* 82:131–142
- Mahmood K, Batool SA, Chaudhery MN, Ul-Haq Z (2017b) Ranking criteria for assessment of municipal solid waste dumping sites. *Arch Environ Prot* 43(1):97–107
- Mahmood K, Ul-Haq Z, Faizi F, Tariq S, Muhammad AN, Rana AD (2019) Monitoring open dumping of municipal waste in Gujranwala, Pakistan using a combination of satellite based bio-indicators and GIS analysis. *Ecol Ind* 107:105613
- Maina CW, Raude JM (2016) Assessing land suitability for rainwater harvesting using geospatial techniques: a case study of Njoro catchment, Kenya. *Appl Environ Soil Sci* 2016:1–9
- Manzo C, Mei A, Zampetti E, Bassani C, Paciucci L, Manetti P (2017) Top-down approach from satellite to terrestrial rover application for environmental monitoring of landfills. *Sci Total Environ* 584–585:1333–1348
- Mugo GM, Odera PA (2019) Site selection for rainwater harvesting structures in Kiambu County-Kenya. *Egypt J Remote Sens Space Sci* 22:155–164
- Nabi G, Ali M, Khan S, Kumar S (2019) The crisis of water shortage and pollution in Pakistan: risk to public health, biodiversity, and ecosystem. *Environ Sci Pollut Res* 26(11):10443–10445
- Ponce VM, Hawkins RH (1996) Runoff curve number: has it reached maturity? *J Hydrol Eng* 1996:11–19
- Satheeshkumar S, Venkateswaran S, Kannan R (2017) Rainfall–runoff estimation using SCS–CN and GIS approach in the Pappiredipatti watershed of the Vaniyar sub basin, South India. *Model Earth Syst Environ* 3:24
- Sekar I, Randhir TO (2007) Spatial assessment of conjunctive water harvesting potential in watershed systems. *J Hydrol* 334(1–2):39–52
- Shukur HK (2017) Estimation curve numbers using GIS and Hec-GeoHMS model. *J Eng* 23(5):1–11
- Soulis KX, Valiantzas JD, Dercas N, Londra PA (2009) Investigation of the direct runoff generation mechanism for the analysis of the SCS-CN method applicability to a partial area experimental watershed. *Hydrol Earth Syst Sci* 13:605–615
- Tumbo SD, Mbilinyi BP, Mahoo HF, Mkilamwinyi FO (2012) Identification of suitable indices for identification of potential sites for rainwater harvesting. *Tanzan J Agric Sci* 12(2):35–46
- USDA (1974) Soil classification system. Definition and abbreviations for soil description. West technical service center, Portland, Oregon, USA
- Yan WY, Mahendrarajah P, Shaker A, Faisal F, Luong R, Al-Ahmad M (2014) Analysis of multi-temporal Landsat satellite images for monitoring land surface temperature of municipal solid waste disposal sites. *Environ Monit Assess* 186(12):1861–1873



# A note on the potential impact of aviation emissions on jet stream propagation over the northern hemisphere

Magdalena Kossakowska<sup>1</sup> · Jacek W. Kaminski<sup>1,2</sup>

Received: 22 August 2019 / Accepted: 8 May 2020 / Published online: 29 May 2020  
© The Author(s) 2020

## Abstract

The goal of the study was to investigate if aviation emissions could influence the climate and weather by modifying the chemical composition of the atmosphere and subsequently, the radiative balance. To carry out the set objective, we used the global environmental multiscale atmospheric chemistry model with comprehensive tropospheric and stratospheric chemistry that is interactive with the radiation calculations. The model was run for two current climate scenarios, with and without aviation emissions. The results of the study indicate that the most significant difference in the jet stream propagation occurred during the winter season, and the smallest was observed during summer. Changes in the jet stream propagation vary by season and region. During the colder time of the year, the eddy-driven jet stream tends to shift poleward, while during the spring season the equatorward shift was observed in a scenario with aviation emissions. Analysis of regional changes shows that the most noticeable differences occurred over the Pacific Ocean, Atlantic Ocean and Asia. The changes over the oceans changed the occurrence of the North Pacific and Bermuda–Azores Highs. Over Asia (Siberia), a stronger and more poleward drift of the eddy-driven jet stream was observed in a scenario without aviation emission. Dissimilarity in the jet stream velocity was found only during the winter seasons when in a scenario with aviation emission, the jet stream velocity was 10 m/s smaller as compared to the scenario without aviation emission.

**Keywords** Aviation emissions · Jet stream · Upper troposphere lower stratosphere · Global environmental multiscale atmospheric chemistry model (GEM-AC)

## Introduction

One of the most significant regions of the atmosphere is the tropopause layer, called the upper troposphere and the lower stratosphere region (UTLS). UTLS is a transition layer where the boundary between the polluted troposphere and ozone-rich stratosphere lies. It plays an important role in tropospheric large scale circulation, stratosphere-troposphere exchange (STE) and the quasi-biennial oscillation (QBO) in the stratosphere (i.e. Holton 1995; Jensen et al. 1996; Garfinkel and Hartmann 2010; Forster and Shine 1997). Any changes in the chemical composition of this region will lead to changes in the dynamics through changes

in the radiative processes (Brasseur et al. 2008; Gettelmann et al. 2011; Hegglin et al. 2010; Shepherd 2002, 2007).

Anthropogenic pollution has a significant impact on atmospheric composition in the troposphere. Most of the sources are near the ground. Thus, the majority of the chemical reactions will take place in the lower and the middle troposphere. Only inert and a small number of reactive species from the ground-based anthropogenic emissions reach the upper troposphere. The aviation emissions, on the other hand, are released mostly in the UTLS region (Olsen et al. 2013a, b). That may cause significant changes in the atmospheric chemistry near the tropopause, especially in the area of heavy airline traffic. Analysis of different aircrafts' fuel burn datasets indicates that 69.0% of aviation emissions are released over the mid-latitudes of the Northern Hemisphere (with the maximum at 40° N), especially over North America (with maximum at 90° W), Europe (maximum between 0° and 10° E) and East Asia (maximum over 115° E). Almost 75% of aircrafts' fuel

✉ Magdalena Kossakowska  
mkossakowska@igf.edu.pl

<sup>1</sup> Institute of Geophysics, Polish Academy of Sciences,  
Ksiecicia Janusza 64 Street, 01-452 Warsaw, Poland

<sup>2</sup> WxPrime Corporation, 21 St. Clair Ave East, Suite 1005,  
Toronto, ON M4T 1L9, Canada

burn takes place in the UTLS, at the height of 7 km (Wilkerson et al. 2010; Olsen et al. 2013a, b).

The complex interactions of gaseous species, direct and indirect effects of aerosols, as well as aviation contrails on the atmospheric chemistry and microphysics, make it difficult to estimate the potential impact of aviation emissions on climate (Penner 1999 (IPCC); IPCC AR5 2014). Lee et al. (2009, 2010) estimated that for the year 2005, aviation was responsible for about 3.5% of the total anthropogenic radiative forcing, including aviation induced cloudiness. This contribution increases up to 4.9% with a range of 2% to 14% for a 90% likelihood range. Many studies show how sensitive modelling results are to aviation emissions and their changes in the UTLS region. For example, lowering flight altitude would lead to changes in the radiative forcing near the tropopause due to the increase in the upper troposphere's ozone mixing ratio. Increasing the flight altitude would lead to the injection of aviation emissions directly into the stratosphere that may have a significant influence on radiative processes (Frömming et al. 2012; Jacobson et al. 2012; Skowron et al. 2015; Søvde et al. 2014).

Most studies focusing on the impact of aviation emissions on climate calculate the global or regional climate change indicator like the mean temperature, radiative forcing or GWP100. However, available studies do not show the exact influence of aviation emissions on global circulation. In the presented study, we decided to examine the sensitivity of the jet stream propagation to aviation emissions, as an indicator of changes in global circulation. We can assume that due to changes in temperature over the Arctic (Yang et al. 2019; IPCC AR5, 2014; Jacobson et al. 2012) or in the low latitude upper troposphere (Grewe et al. 2002; Lee et al. 2010; Lund et al. 2017) there may be a noticeable change in the jet stream propagation that can strongly affect some regions, especially over the mid- and high latitudes (Barnes and Simpson 2017; Cohen et al. 2014; Linz et al. 2018; Xue et al. 2017). Studies suggest the general poleward shift of the eddy-driven jet stream (EDJ) as well as the subtropical jet stream (STJ), but those trends vary, depending on seasons or regions (Melamed-Turkosh et al. 2018; Rikus 2018; Strong and Davis 2007; Zolotov et al. 2018).

The tendency of the jet stream to poleward or equatorward shifts is mostly driven by the upper troposphere tropical warming (Sun et al. 2013; Simpson et al. 2012) and Arctic warming (Barnes and Simpson 2017), respectively. The Arctic warming will slow down the poleward jet stream shift due to GHG impact on low latitudes (Barnes and Polvani, 2013; Barners and Screen, 2015; Haigh et al. 2005; Linz et al. 2018). On the regional scale, the changes in the jet stream propagation may be influenced by sea surface temperature, ice cover, ENSO, stratospheric polar vortex, radiative forcing, QBO or volcanic eruptions (Hall et al. 2015).

## Method

The objective of the presented study was to examine changes in jet stream propagation due to aviation emissions. We designed two current climate modelling scenarios: base scenario A0 without aviation emissions and scenario A1 with aviation emissions. We used the Global Environment Multiscale model GEM-AC with interactive and coupled tropospheric and stratospheric chemistry (de Grandpré et al. 2000; Kaminski et al. 2008; Mamun et al. 2013; Lupu et al. 2013). The model horizontal grid was defined as the global variable resolution from  $3^\circ \times 3^\circ$  to  $1.5^\circ \times 1.5^\circ$  zoomed over the high latitudes of the Northern Hemisphere, starting from  $55^\circ$  N (Fig. 1) with 70 hybrid vertical levels up to 0.1 hPa and a 30 min time step and output set at every 6 h. The vertical resolution in the UTLS region was 500 m.

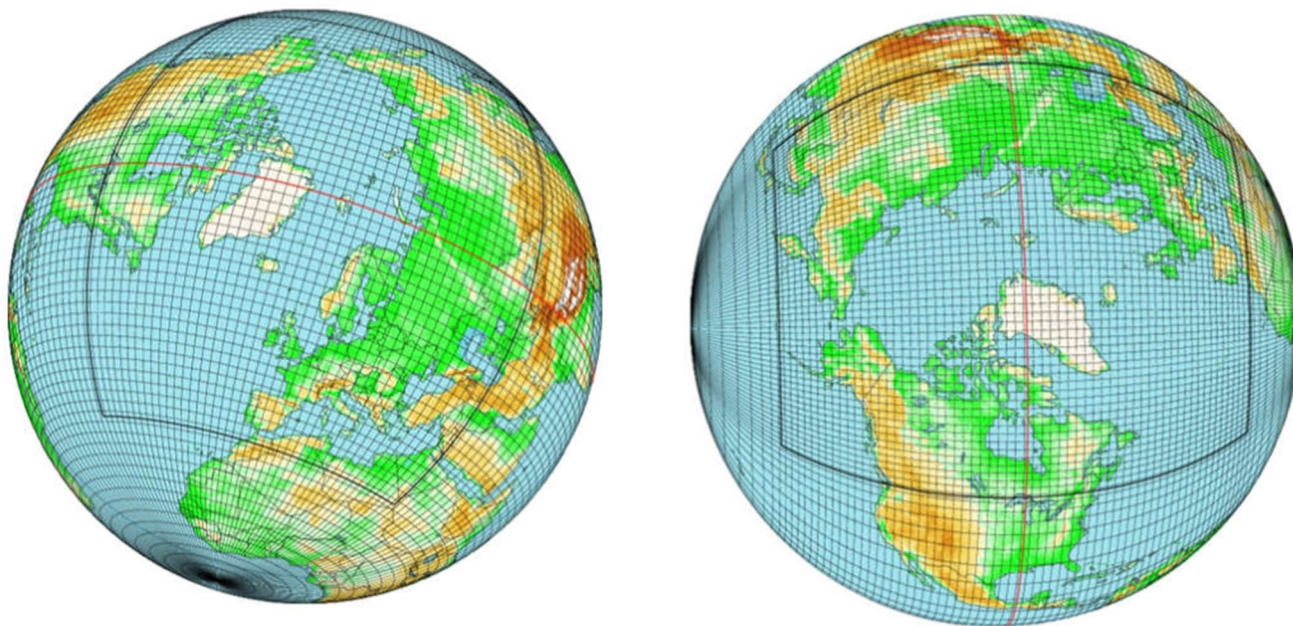
Both simulations were run in a climate mode setup for years 2001–2010. Both scenario runs were started with a “cold” model, allowing chemical species to balance in the atmosphere during the first 5 years of the simulation. Due to the model “cold” run, only results after 2005 could be analysed. Inventories of aviation emissions were available only for 2006, while for other years the aviation emissions were rescaled based on estimations provided by the International Civil Aviation Organisation (ICAO). In this paper, only results for the year 2006 are presented.

Climatological information is based on monthly mean ice cover and sea surface temperature, obtained from the geophysical fluid dynamics laboratory (GFDL) model. Historical anthropogenic emissions (excluded aviation) were taken from ACCMIP (Lamarque et al. 2012). Aviation emissions used in scenario A1 were from AEDT 2006 database provided as hourly 3D fields of the total fuel burn and CO, HC, NO<sub>x</sub>, PM<sub>NV</sub>, PM<sub>SO</sub>, PM<sub>FO</sub>, CO<sub>2</sub>, H<sub>2</sub>O, SO<sub>x</sub> with horizontal of  $1^\circ \times 1^\circ$  and 500 ft in vertical (Kim et al. 2007; Wilkerson et al. 2010). The initial conditions for GEM-AC were generated using CMAM (de Grandpre et al. 2000, 2009).

For this study, the jet stream was defined as a narrow, horizontal air current with a wind speed greater than 25 m/s, located between 400 and 100 hPa. The definition is based on WMO's jet stream description (1958), and simple jet stream detection method proposed by Pena-Ortic et al. (2013) was used.

## Results

Results for 2006 simulations with (scenario A1) and without (scenario A0) aviation emissions were analysed using the annual, seasonal (winter (DJF), spring (MAM),

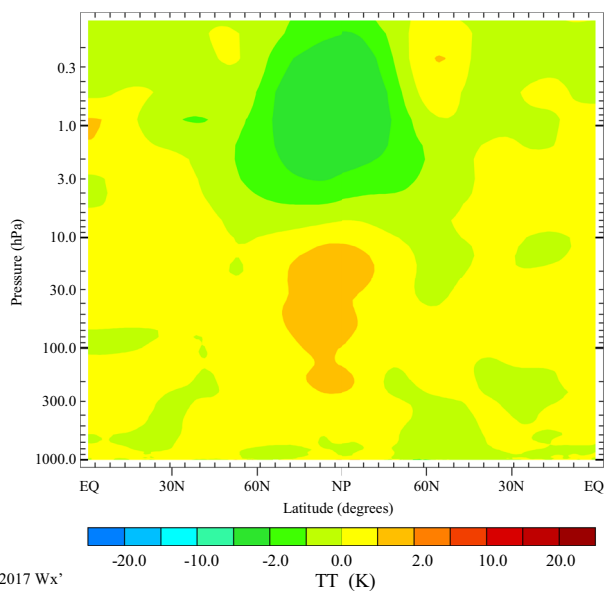


**Fig. 1** Visualisation of the GEM-AC model grid with global variable horizontal resolution  $3.0^\circ \times 3.0^\circ$  and  $1.5^\circ \times 1.5^\circ$  regional nested over the Northern Hemisphere. European (left) and the North American (right) vantage point of view

summer (JJA), autumn (SON)) and monthly time interval averaging of the meridional and zonal wind velocities. For each period, we created three different visualizations:

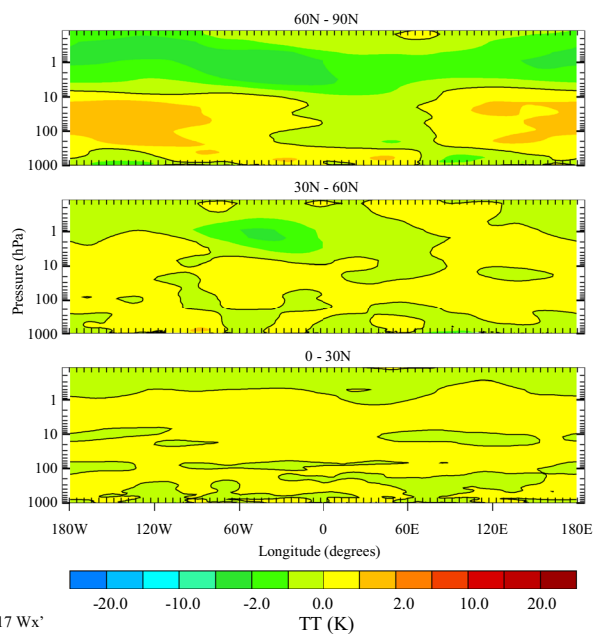
average by longitudes, separately for Western and Eastern half of the Northern Hemisphere (shown in Fig. 2), average by latitudes in three bands: low (0–30), middle

W & E hemispheric average for A1-A0\_3x3\_2006\_YEAR



(c) 2017 Wx'

Average latitude bands for A1-A0\_3x3\_2006\_YEAR



(c) 2017 Wx'

**Fig. 2** Left: The difference between annual mean temperature between scenarios A1 and A0, averaged over longitudes, separately for Western (W) and Eastern half of the Northern Hemisphere. Right:

The difference between annual mean temperature scenarios A1 and A0, averaged over latitudes for low (0–30 N), mid (30–60 N) and high (60–90 N) latitudes bands

(30–60) and high (60–90) latitudes (shown in Fig. 2) and standard zonal average for the whole globe. The jet streams propagation was compared using monthly averaged wind velocity at several isobaric levels. Results for 250 hPa level are shown as both the subtropical and eddy-driven jet streams. Also, we analysed changes in the temperature using analogical methods for the wind velocity to focus on changes of temperature, especially between high latitudes and low latitudes. We focused on changes in temperature near the ground over the polar regions, polar UTLS and the tropical upper troposphere transition layer (TTL).

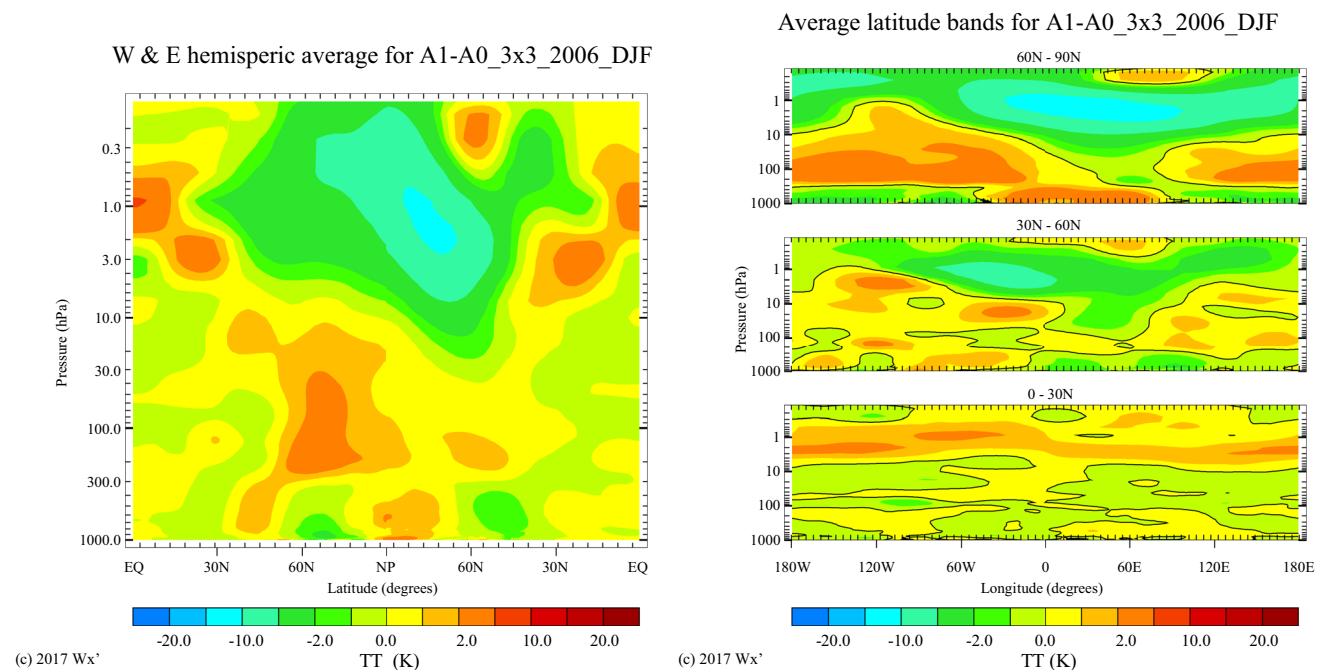
### Changes in the temperature

Analysing the differences between A1 and A0 scenarios' annual mean temperatures, we noticed that the largest changes occurred in the upper troposphere lower stratosphere region over the high latitudes, where scenario A1 shows a temperature that was 2 K higher (Fig. 2). On seasonal and monthly time scales, changes between scenarios were much larger. In the winter season, the influence of aviation emissions on temperature was the strongest over the lower troposphere of high latitudes, especially over the Arctic and Scandinavia, where the seasonal mean temperature in scenario A1 was 4 K higher than in scenario A0 (shown in Fig. 3), and the monthly differences vary from 3 K up to 8 K. Also, we noticed regions with significant negative

temperature changes in the high latitudes, but with strong regional variation from month to month where changes reached  $-5$  K.

During the winter season, the aviation emissions lead to a small temperature decrease in the low latitudes' upper troposphere indicate (up to  $-2$  K). The changes in the mid-latitudes' troposphere vary, depending on region, from  $-2$  K (mainly over the Eastern NH) up to 2 K (mostly over the Western NH), as can be seen on Fig. 3. In the mid-latitude UTLS region, the aviation emissions mainly cause small (up to  $-2$  K) temperature increase. The changes due to aviation emissions in the high latitudes' lower troposphere indicate a  $-5$  K decrease over the region between mid- and high latitudes but up to 5 K increase over the Arctic region. In the UTLS, we were able to observe up to 5 K temperature increase.

The analysis of the monthly regional changes shows pronounced temperature differences between scenarios, especially over the high latitudes. The largest differences between scenarios occurred over Russia, Europe, Baffin Bay, the Bering Sea and Central North America, with temperature variation from  $-9$  K up to 10 K, from  $-6$  K to 0 K, from  $-9$  K to 5 K, from  $-10$  K to 5 K and from  $-7$  K up to 3 K, respectively. Only over the central and eastern part of the USA, there was a constant increase in temperature in scenario A1, with the differences between scenarios of less than 4 K. In the UTLS region, we noticed a small decrease



**Fig. 3** Left: The difference between the seasonal mean temperature between scenario A1 and A0, averaged over longitudes, separately for Western (W) and Eastern half of the Northern Hemisphere for the winter season. Right: The difference between the seasonal mean tem-

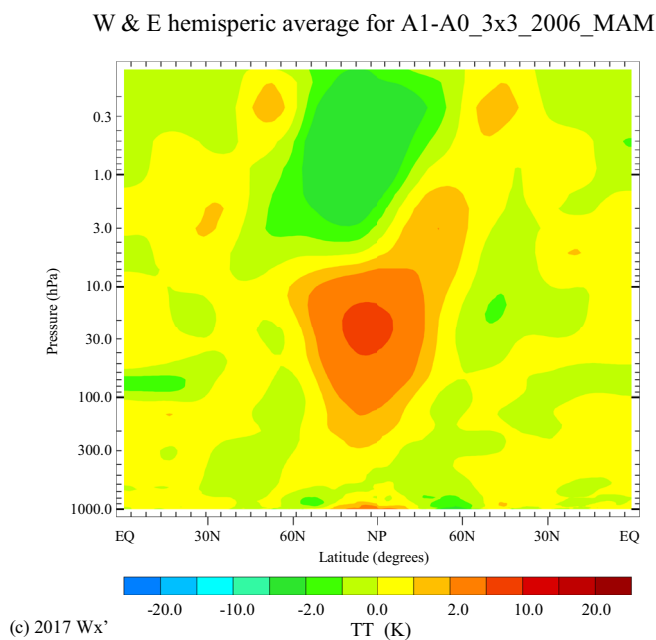
perature between scenario A1 and A0, averaged over latitudes for low (0–30 N), mid (30–60 N) and high (60–90 N) latitudes bands for the winter season

(−3 K) over the Atlantic Ocean between 30° N and 45° N. A similar tendency occurred over the Pacific Ocean for the same latitude range, but with stronger variation along the longitudes. The increase in the UTLS temperature occurred mostly over the Western US, Canada, the Labrador Sea and Greenland, with the highest difference of up to 5 K. Also, a small temperature increase was noticed over the eastern part of Russia, where the differences between scenario A1 and A0 did not exceed 3 K. Over Europe, monthly differences varied from 3 K to −4 K.

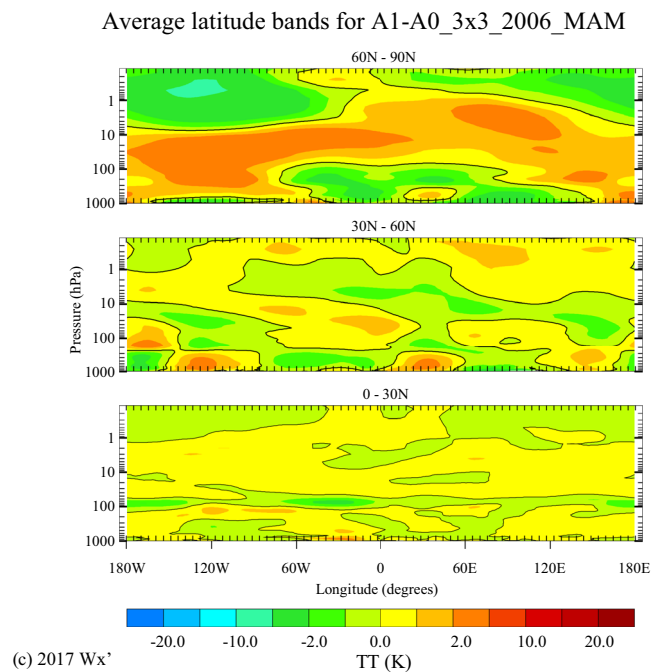
During the spring season, the analysis of changes between scenarios A1 and A0 also indicates a small (up to −2 K) temperature decrease in the low latitudes upper troposphere due to aviation emissions. In the mid-latitudes, the influence of aviation emissions on the troposphere temperature leads to substantial regional variation, mostly related to the ocean–land presence (please see Fig. 4, right panel). Over the region with land domination (except Eastern Europe), we noticed up to 5 K temperature increase in the troposphere. In contrast, for the region with ocean domination—scenario A1 shows up to −5 K lower mean meridional temperature. That may be connected with the stronger aviation emissions in the lower and mid-troposphere over the land, because of the presence of airports, while over the oceans aviation emissions take place at cruising altitudes because of the small number of airports at isolated archipelagos. In the mid-latitudes UTLS region, the temperature changes due to

aviation emissions were opposite to the trend we noticed in the troposphere and varies from −2 K over the lands up to 2 K over the oceans. The mean zonal temperature difference for spring season shows the small west–east hemispheric contrast with the prevailing warming effects in the Eastern NH and cooling effect in the Western NH of the aviation emissions in the troposphere. Over the high latitudes, we noticed the same general trends in temperature changes as we did for winter months. The temperature increases (up to 2 K) over the Arctic, but there is a small (−2 K) decrease around the Arctic between 60° N and 70° N. In the high latitudes UTLS region, the aviation emissions scenario shows up to 5 K temperature increase (especially in the lower stratosphere).

For the summer, the decreasing trend in the low latitude upper troposphere temperature due to aviation emissions visible during the winter and spring seasons start to change, showing the regional upper troposphere temperature increase. The analysis of monthly zonal temperature means for low latitudes indicate month to month trend changing from temperature decrease to temperature increase in scenario A1. In the mid-latitudes troposphere, the seasonal mean temperature difference between scenarios was rather small and oscillated between  $\pm 1$  K. The changes in the mid-latitudes UTLS indicate a similar regional variation as in the troposphere, but with the opposite sign of changes (see Fig. 5, right middle panel). These regional changes follow

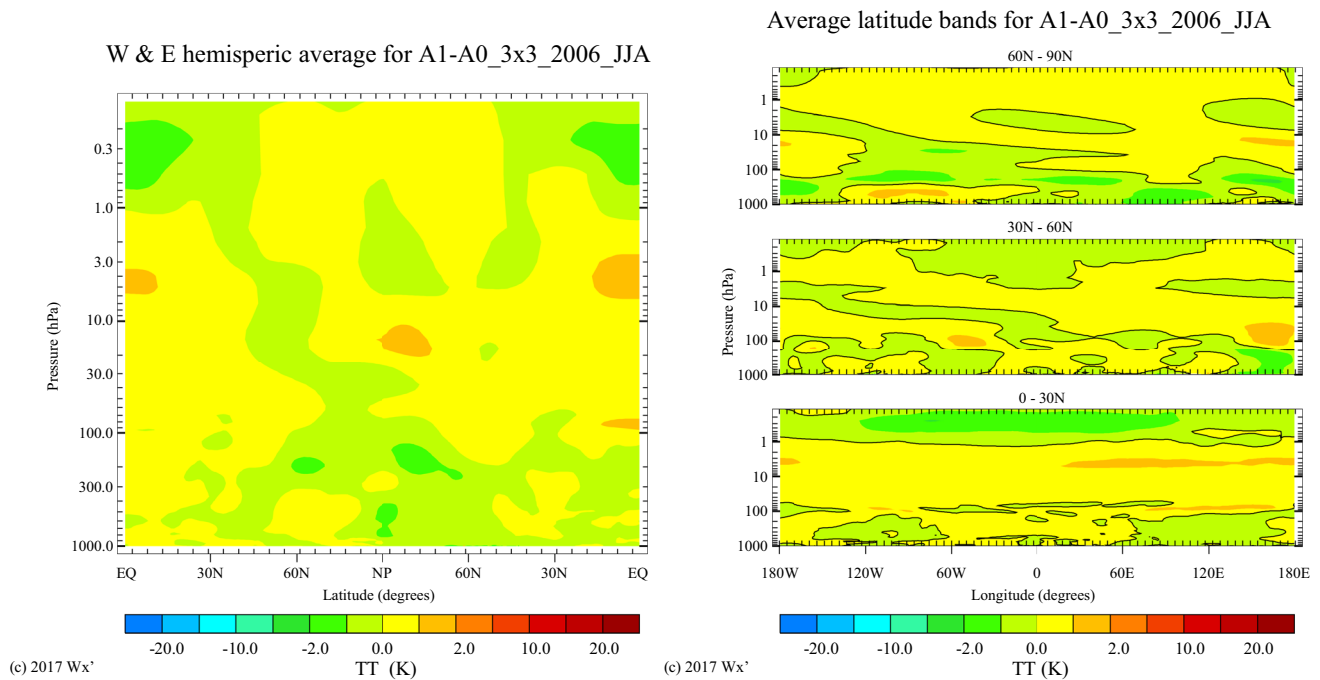


**Fig. 4** Left: The difference between the seasonal mean of temperature between scenario A1 and A0, averaged over longitudes, separately for Western (W) and Eastern half of the Northern Hemisphere for the spring season. Right: The difference between the seasonal mean of



temperature between scenario A1 and A0, averaged over latitudes for low (0–30 N), mid (30–60 N) and high (60–90 N) latitudes bands for the spring season





**Fig. 5** Left: The difference between the seasonal means of temperature between scenario A1 and A0, averaged over longitudes, separately for Western (W) and Eastern half of the Northern Hemisphere the summer season. Right: The difference between the seasonal

means of temperature between scenario A1 and A0, averaged over latitudes for low (0–30 N), mid (30–60 N) and high (60–90 N) latitudes bands for the summer season

the trend described for the spring season. In the high latitudes, we noticed a decreasing trend in temperature due to aviation emissions. The exception occurred between 30° W and 120° W, over the northern part of Canada and Greenland, where we noticed up to 5 K increase.

During the autumn season, the direction of changes in the low latitudes upper troposphere indicates a 2 K increase in temperature due to aviation emissions. The changes in the troposphere and UTLS mean temperature over the mid-latitudes have a robust regional variability, similar to the trend described for spring and summer season with a small spatial shift. The overall mean seasonal changes indicate a small (1 K) tropospheric temperature increase over the Western NH and a small (–1 K) temperature decrease over the Eastern NH. Over the high latitudes, the changes in the lower troposphere indicate a –2 K decrease over the Arctic with 1 K increase over Northern Canada and Scandinavia. The changes in the high latitude UTLS vary from –2 K over the Siberia region up to 5 K over North Canada. The mean seasonal changes over the high latitudes UTLS during autumn indicate a 2 K increase in scenario A1 (Fig. 6).

The mean annual changes between scenarios presented in Fig. 2 show no noticeable changes between scenarios. However, the analysis shows significant monthly and regional differences between scenarios A1 and A0. Although the most intense aviation emissions occur over the mid-latitudes,

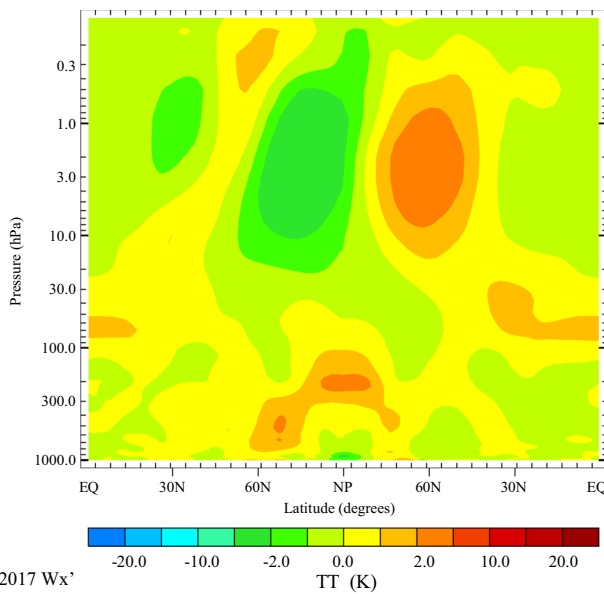
the most sensitive regions seem to be high and near high latitudes, where the changes vary in a range of  $\pm 7$  K. The changes in the temperature over the mid-latitudes show a less strong response to aviation emissions, yet the changes strongly depend on the emission area. When aviation emissions were present in the whole troposphere, we observed the temperature increase in the low and middle troposphere but decrease in the UTLS region. On the one hand, when aviation emissions were limited only to cruise altitudes temperature in the UTLS increase while in the middle and low troposphere, we noticed a temperature decrease. The influence of aviation emissions on the tropical upper troposphere's temperature was the weakest, yet we still noticed some small changes.

### Changes in the jet stream propagation

The analysis of annual means shows almost no differences between scenarios A1 and A0 for the zonal wind velocity and only small shifts in the meridional jet stream velocity over the North America mid-latitudes (Fig. 7). However, differences between scenarios in seasonal and monthly mean values show noticeable changes in the jet stream propagation.

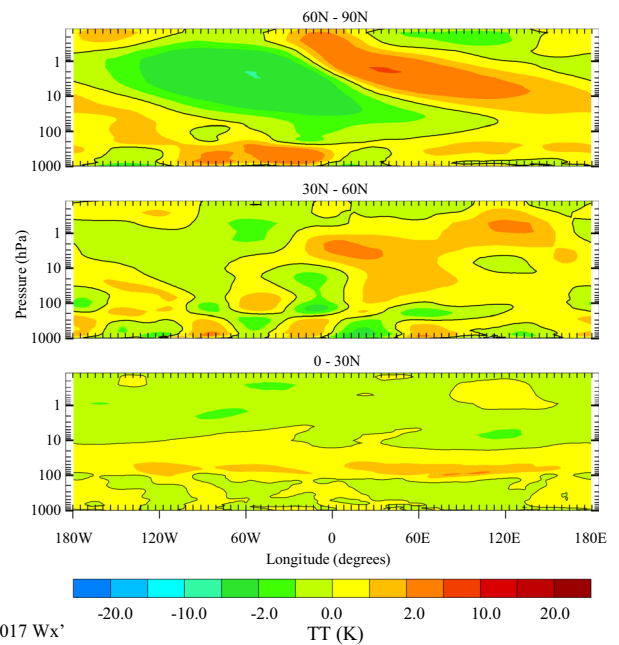
The strongest influence of aircraft emissions on jet stream was noticed during the winter season, where

W & E hemispheric average for A1-A0\_3x3\_2006\_SON



(c) 2017 Wx'

Average latitude bands for A1-A0\_3x3\_2006\_SON

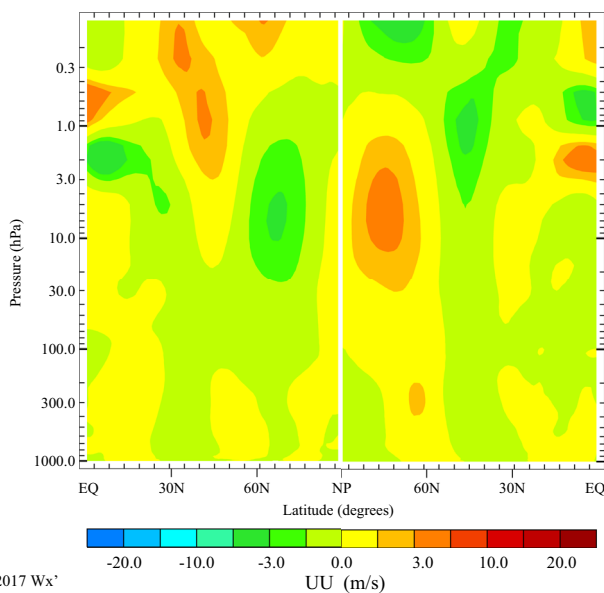


(c) 2017 Wx'

**Fig. 6** Left: The difference between the seasonal means of temperature between scenario A1 and A0, averaged over longitudes, separately for Western (W) and Eastern half of the Northern Hemisphere for the autumn season. Right: The difference between the seasonal

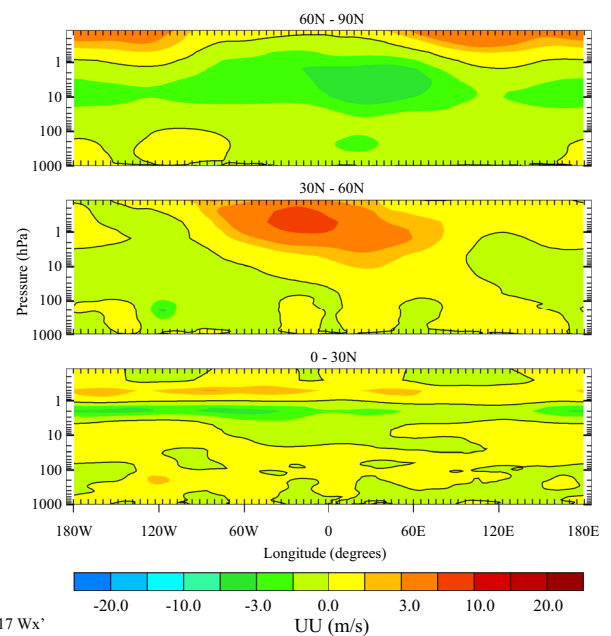
means of temperature between scenario A1 and A0, averaged over latitudes for low (0–30 N), mid (30–60 N) and high (60–90 N) latitudes bands for the autumn season

W & E hemispheric average for A1-A0\_3x3\_2006\_YEAR



(c) 2017 Wx'

Average latitude bands for A1-A0\_3x3\_2006\_YEAR



(c) 2017 Wx'

**Fig. 7** Left: The difference between annual means of zonal wind velocity between scenario A1 and A0, averaged by longitudes, separately for Western (W) and Eastern half of the Northern Hemisphere.

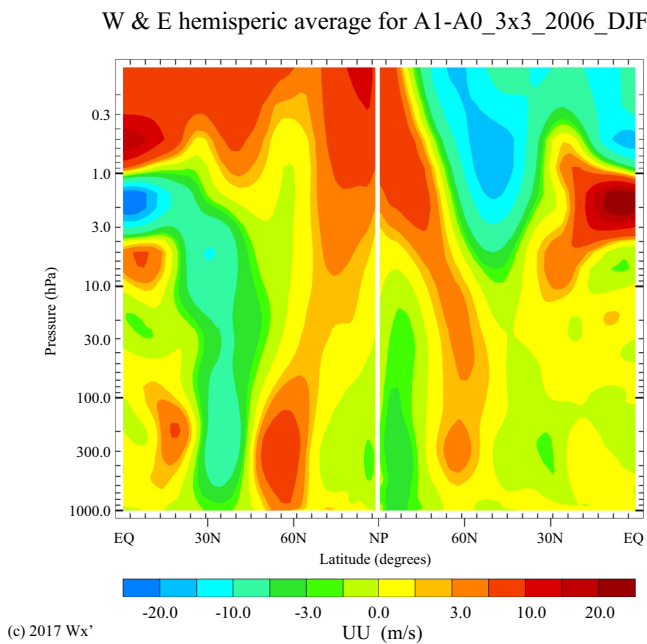
Right: The difference between annual means of meridional wind velocity between scenario A1 and A0, averaged over latitudes for low (0–30 N), mid (30–60 N) and high (60–90 N) latitudes bands

aviation caused an increase in zonal wind velocity in the UTLS region over the high and low latitudes and a decrease in mid-latitudes, indicating the separation of the STJ and EDJ. The poleward shift of the polar jet stream and equatorward shift of the STJ was observed mostly over the Western Hemisphere. The changes over the Eastern Hemisphere in the seasonal wind speed analysis are less visible. Analysis of the regional changes in the jet stream propagation shows that the largest differences between scenarios occurred over Northern America, the Atlantic Ocean, East Asia and the Pacific Ocean. In the scenario with aviation emissions, the jet stream has a tendency to split into two streams and propagate on both sides of the Rocky Mountains, travelling more often over Northern Canada and more often over the southern part of the USA. Those stronger tendencies to propagate along the Rocky Mountains can be seen in changes in the meridional wind velocity, shown in Fig. 8. Both the EDJ and STJ jets stay separated over the Atlantic Ocean. The tendency of the jet stream to split over North America results in poleward (EDJ) and equatorward (STJ) shifts we noticed in the seasonal mean of the wind velocity in the UTLS region. Also, we noticed a tendency of jet stream superposition over Asia, which cause a decrease in EDJ occurrence over the Eastern Siberia and wind velocity increase over East Asia

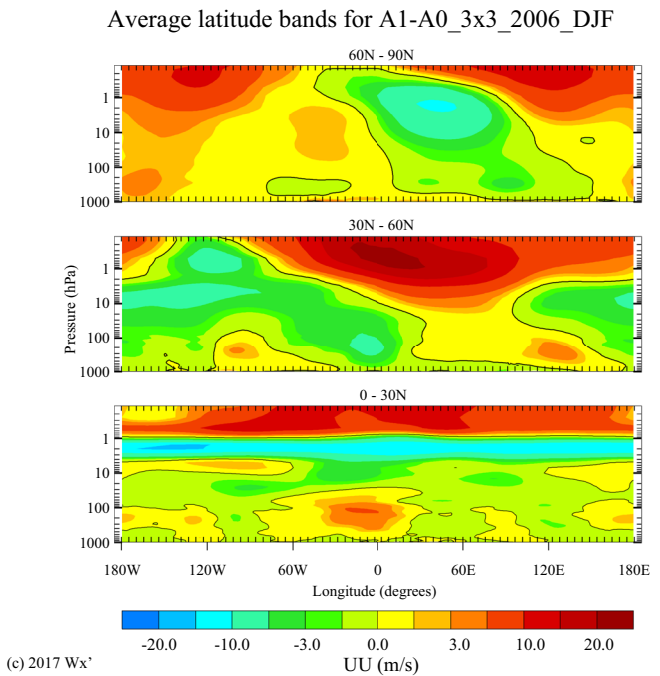
(China). The tendency for January is shown as an example in Fig. 9.

The analysis of differences in the jet core velocity between A1 and A0 scenarios shows small, not noticeable changes in the seasonal mean of the jet core velocity. However, the regional differences between scenarios can reach up to  $-10$  m/s. The largest difference occurred over the Pacific Ocean in January, when the monthly mean of the jet core velocity in scenario A1 was around 66 m/s while in scenario A0 it was 77 m/s. During the spring season, the changes between scenarios were less visible as compared to the winter season. The direction of changes in the spring season did not follow the winter trend. The jet streams in the scenario with aviation emissions show a small equatorward shift with a stronger and more stable subtropical jet. The changes in the jet stream propagation vary between scenarios on a month to month time scale. Still, the mean seasonal trend shows that stronger changes occurred over the Western Northern Hemisphere (Fig. 10).

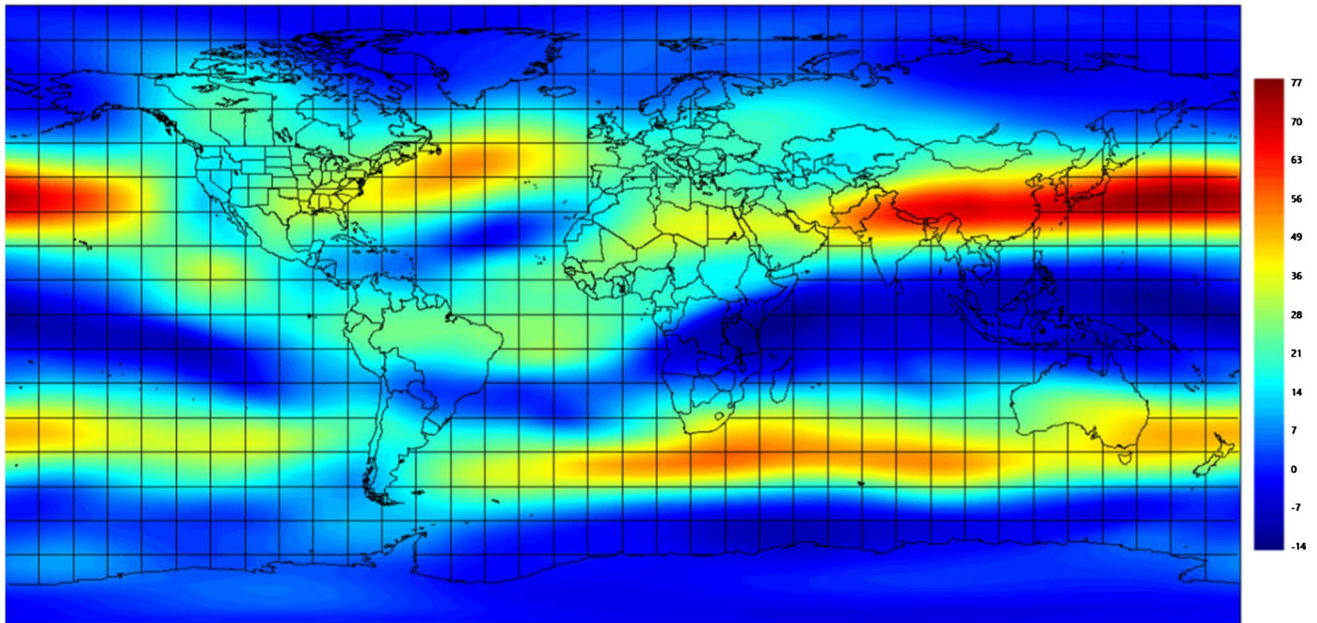
The regional changes in the jet streams propagation due to aviation emissions mostly affected the Pacific and Atlantic Oceans. Analysis of the monthly means of wind velocity field indicates that in scenario A1 the jet stream tends to split over the East Pacific, where EDJ propagate poleward toward Alaska and Canada, and STJ propagate equatorward toward Mexico. This tendency follows the trend from the winter



**Fig. 8** Left: The difference between and seasonal means of zonal wind velocity between scenario A1 and A0, averaged over longitudes, separately for Western (W) and Eastern half of the Northern Hemisphere for the winter season. Right: The difference between seasonal

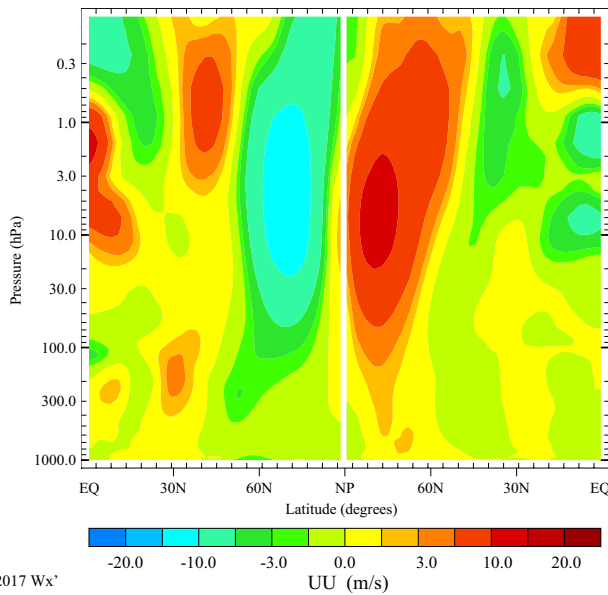


means of meridional wind velocity between scenario A1 and A0, averaged over latitudes for low (0–30 N), mid (30–60 N) and high (60–90 N) latitudes bands for the winter season



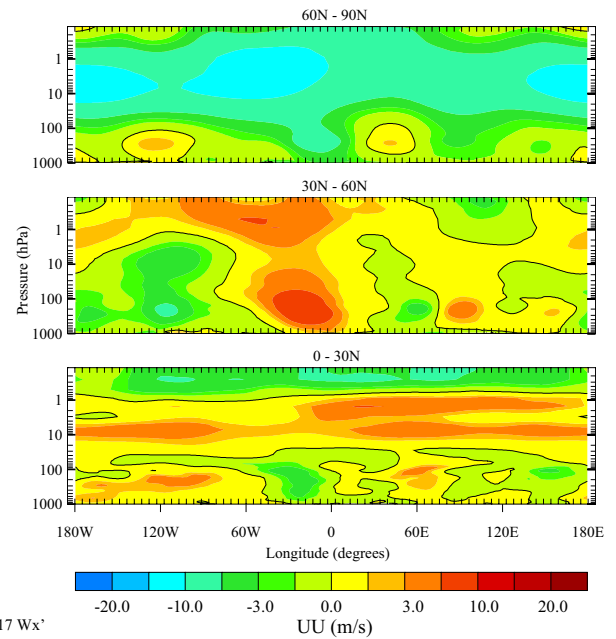
**Fig. 9** Monthly mean wind velocity at 250 hPa in knots for scenario A1 (top panel) and scenario A0 middle panel). The bottom panel shows a difference between monthly mean wind velocity between scenarios A1 and A0 at 250 hPa in knots. 1 knot is equal to 0.514 m/s

W & E hemispheric average for A1-A0\_3x3\_2006\_MAM



(c) 2017 Wx'

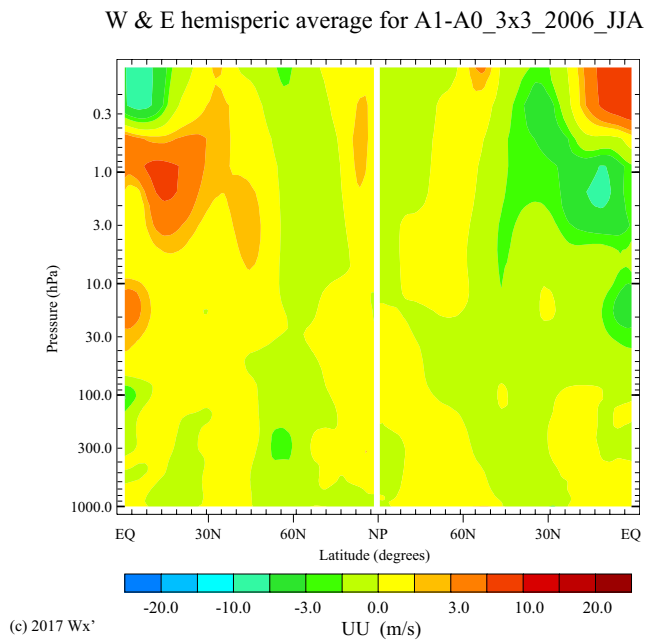
Average latitude bands for A1-A0\_3x3\_2006\_MAM



(c) 2017 Wx'

**Fig. 10** Left: The difference between seasonal means of zonal wind velocity between scenario A1 and A0, averaged over longitudes, separately for Western (W) and Eastern half of the Northern Hemisphere for the spring season. Right: The difference between seasonal means

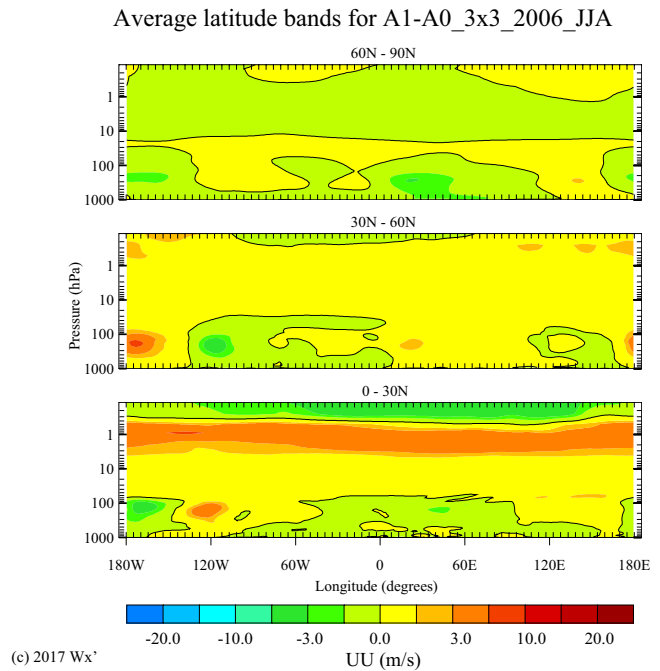
of meridional wind velocity between scenario A1 and A0, averaged over latitudes for low (0–30 N), mid (30–60 N) and high (60–90 N) latitudes bands for the spring season



**Fig. 11** Left: The difference between seasonal means of zonal wind velocity between scenario A1 and A0, averaged over longitudes, separately for Western (W) and Eastern half of the Northern Hemisphere for the summer season. Right: The difference between seasonal

season. In scenario A0, this tendency is much weaker and appears at the end of the spring season. On the other hand, in the scenario without aviation emissions, we noticed a strong tendency to form the North Atlantic Subtropical High that split the jet flow into two streams. The blocking pattern of the Bermuda-Azores high is not visible in the results for the scenario A1.

The smallest changes in the jet stream propagation due to the aviation emissions were found during the summer season (Fig. 11). Analysis of the monthly zonal means of the wind velocity has shown the small variation in the zonal jet stream flow between scenarios. Small shifts in the jet stream propagation, visible in monthly means of the zonal wind speed in the upper troposphere, show no particular trend in the jet stream modification between June and August for A1 and A0 scenarios. For the meridional wind velocity, we noticed meander-like structures in the jet stream propagation over the mid-latitudes. The monthly means show that in scenario A1 the jet stream tends to wobble more often over North America. In contrast, in scenario A0, the meanders are more frequent over Asia and the Pacific Ocean. The changes over North America were connected to the stronger trend of the jet stream split over the East Pacific Ocean that we noticed in scenario A1 from the winter season. The differences over Asia and the West Pacific Ocean were due to an increase in the tendency to jet stream split over Asia in the scenario with aviation emissions. This tendency results in more frequent

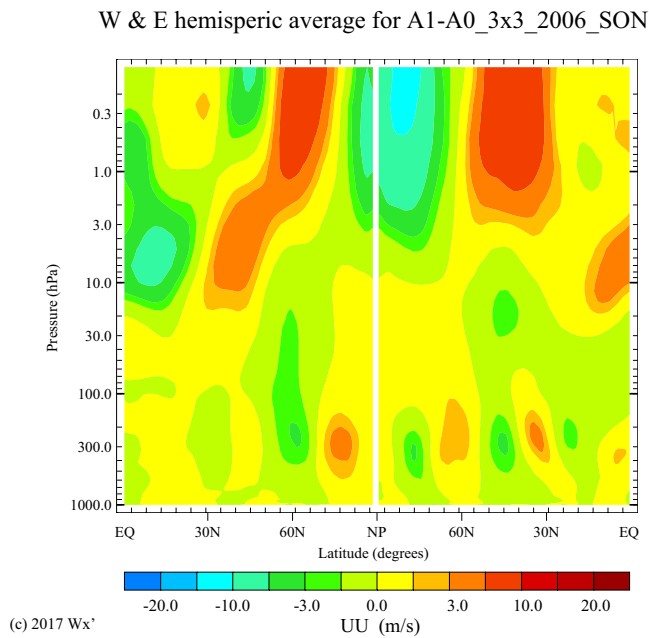


means of meridional wind velocity between scenario A1 and A0, averaged over latitudes for low (0–30 N), mid (30–60 N) and high (60–90 N) latitude bands for the summer season

EDJ occurrence over the high latitudes and equatorward shift of the STJ over the West Pacific Ocean during the summer season.

The analysis of the differences between wind fields for the autumn season indicates a poleward shift of the EDJ and STJ in the scenario with aviation emissions, as presented in Fig. 12. Both scenarios show the tendency to jet stream split over the East Pacific Ocean, but in the scenario A0, this tendency is more noticeable. The differences in meridional flow indicate a significant change in the jet stream propagation, especially over Asia. The analysis of the wind velocity in the UTLS region shows that in both scenarios, the jet streams tend to wobble with similar intensity. There was no noticeable tendency in changes in jet stream meandering between scenarios A1 and A0. The strongest influence of aviation emissions was observed over the Pacific Ocean, Atlantic Ocean and Europe. In the scenario with aviation emissions, the jet stream shows a poleward shift over the Pacific Ocean with a small difference in the mean jet core velocity. On the other hand, over the western part of the Atlantic Ocean, we noticed a significant decrease in jet stream strength and equatorward shift of the jet stream flow over Europe that leads to the increase in the zonal wind velocity over Central and South Europe.

In summary, there is a noticeable month to month and season to season variation in the jet streams propagation between scenarios A1 and A0. The most significant

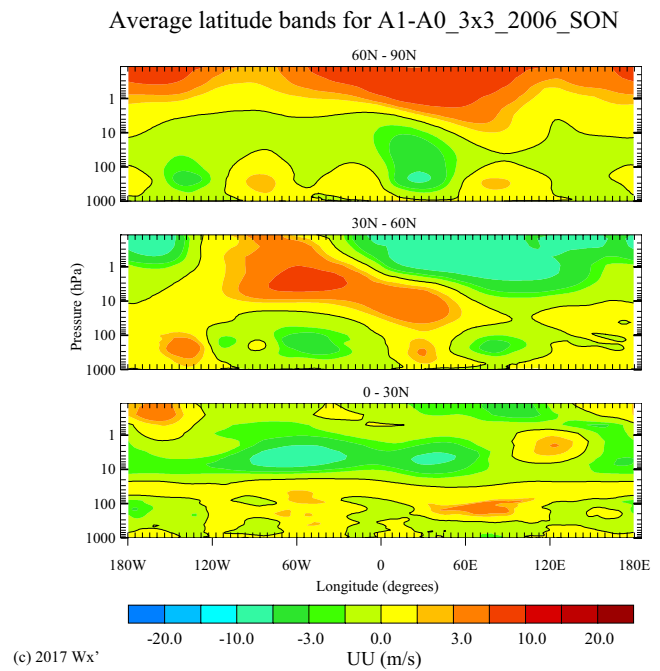


**Fig. 12** Left: The difference between seasonal means of zonal wind velocity between scenario A1 and A0, averaged over longitudes, separately for Western (W) and Eastern half of the Northern Hemisphere for the autumn season. Right: The difference between seasonal means

difference occurred during the winter season and the smallest in summer. There is no regular trend in the jet stream shifts between scenarios. During the colder time of the year, the EDJ tends to shift poleward while during the spring season we noticed an equatorward shift in scenario A1. The analysis of the jet core monthly mean velocity shows no difference between scenarios except for the winter season, when the jet streams in scenario A1 seems to be slower as compared to the results for the A0 scenario. Regional changes in the jet stream propagation due to aviation emissions are mainly visible over the Pacific and Atlantic Oceans, where jet streams tend to be more stable, most likely due to the uniform surface. The most noticeable difference between scenarios is the stronger tendency of the jet stream split in scenario A1, especially over the East Pacific Ocean.

## Summary and conclusions

Changes in the troposphere and the UTLS region temperature due to aviation emissions indicate around 2 K temperature decrease in the tropical upper troposphere region during the winter season, while in the UTLS region over high latitudes we noticed up to 5 K temperature increase. The analysis of changes in the tropopause temperature, presented by Hu and Vallis (2019) for the years 1979–2017 shows that the standard deviation of the mean annual temperature at



of meridional wind velocity between scenario A1 and A0, averaged over latitudes for low (0–30 N), mid (30–60 N) and high (60–90 N) latitudes bands for the autumn season

the tropopause is 1.5 K and 1.0 K, for the low and high latitudes, respectively. We can assume that changes in the tropical tropopause exceed  $1\sigma$  of the mean annual climatological tropopause temperature variation, and for the latitudes the changes due to aviation emissions exceed  $2\sigma$ . It would indicate significant changes in the UTLS temperature, especially over the high latitudes.

The general propagation of the jet stream in scenarios A1 and A0 was in agreement with the results presented by Christenson et al. (2017), Koch et al. (2006), Kuang et al. (2014) or Pena-Ortiz et al. (2013). Also, there were regional differences, especially over North America and North Asia, where scenario A1 shows better agreement with jet stream climatological studies than was expected. The preliminary results of the presented one-year case study show that aviation emissions lead to significant changes in jet stream propagation. Analysis of changes in the jet streams propagation indicates that the aviation emissions lead to more polar and subtropical jet splits than in the scenario without aviation emissions, especially during winter. We noticed a poleward EDJ shift during the colder part of the year that may be caused by changes in the UTLS temperature described in the previous paragraph. During the winter, the warming of the high latitudes UTLS is stronger than the cooling in the tropical upper troposphere that leads to the poleward shift of the jet streams over the Northern Hemisphere in scenario A1 during the winter season. The poleward shift of the jet

streams in the autumn season was mostly connected with warming over the high latitudes rather than changes in the TTL. There were no significant changes in the jet stream velocity except in winter, when the mean seasonal jet core velocity in scenario A1 was about 10 m/s slower than in scenario A0. There was no strong constant tendency of the jet stream more frequent wobble flow, yet there were significant regional differences in the jet streams propagation. The most noticeable changes were observed over the Eastern Pacific, Eastern Atlantic and North-East Asia. Over the Northern Pacific, we observed a significant difference in jet streams split what leads to changes in the North Pacific High (NHP) development. In scenario A1, the persistent high pressure system was much stronger in the spring season and weaker in autumn than in scenario A0. These changes may have a significant influence on drought season, especially over California. Another important change in mid-latitudes was the impact of aviation emissions on Bermuda-Azores High during the spring season, which was caused by changes in the jet stream propagation over the North and Subtropical Atlantic. Over Asia (Siberia), we observed stronger and more poleward EDJ in scenario A0.

Modelling results have shown that aviation emissions alone may have a significant influence on the jet stream propagation that leads to the conclusion that aviation emissions may have a significant influence on the climate. At this point, it is essential to highlight that the presented study covered the early research results, based only on one year of the simulation, focusing on the question “if aviation emissions may influence the jet stream propagation”. The future results will be focused on climatological aspects of changes in the jet stream due to increasing aviation traffic.

**Acknowledgements** MK was supported by the National Science Centre Poland, Grant No. UMO 2013/11/N/ST10/00330 and Grant No. 2016/23/B/ST10/03192. JWK was supported by Transport Canada, Contract No: T8497-120001/01/TOR.

## Compliance with ethical standards

**Conflict of interest** On behalf of all authors, the corresponding author states that there is no conflict of interest.

**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article’s Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article’s Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

## References

- Barnes EA, Polvani L (2013) Response of the Midlatitude jets, and of their variability, to increased greenhouse gases in the CMIP5 models. *J Clim* 26:7117–7135. <https://doi.org/10.1175/JCLI-D-12-00536.1>
- Barnes EA, Screen JA (2015) The impact of Arctic warming on the midlatitude jet-stream: Can it? Has it? Will it? *Wiley Interdiscip Rev Clim Change* 6:277–286. <https://doi.org/10.1002/wcc.337>
- Barnes EA, Simpson IR (2017) Seasonal sensitivity of the Northern Hemisphere jet streams to arctic temperatures on subseasonal time scales. *J Clim* 30:10117–10137. <https://doi.org/10.1175/JCLI-D-17-0299.1>
- Brasseur GP (coordinating lead author) et al (2008) Aviation climate change research initiative: a report on the way forward based on the review of research gaps and priorities, Federal Aviation Administration, available from: [https://www.faa.gov/about/office\\_org/headquarters\\_offices/apl/research/science\\_integrated\\_modeling/accrri/media/ACCRI\\_Report\\_final.pdf](https://www.faa.gov/about/office_org/headquarters_offices/apl/research/science_integrated_modeling/accrri/media/ACCRI_Report_final.pdf). Accessed 22 May 2019
- Christenson CE, Martin JE, Handlos ZJ (2017) A Synoptic climatology of Northern Hemisphere, cold season polar and subtropical jet superposition events. *J Clim* 30:7231–7246. <https://doi.org/10.1175/JCLI-D-16-0565.1>
- Cohen J, Screen JA, Furtado JC et al (2014) Recent arctic amplification and extreme mid-latitude weather. *Nat Geosci* 7:627–637. <https://doi.org/10.1038/ngeo2234>
- de Grandpré J, Beagley SR, Fomichev VI et al (2000) Ozone climatology using interactive chemistry: results from the Canadian middle atmosphere model. *J Geophys Res Atmos* 105:26475–26491. <https://doi.org/10.1029/2000JD900427>
- Forster PMD, Shine KP (1997) Radiative forcing and temperature trends from stratospheric ozone changes. *J Geophys Res* 102:10841. <https://doi.org/10.1029/96JD03510>
- Frömming C, Ponater M, Dahlmann K et al (2012) Aviation-induced radiative forcing and surface temperature change in dependency of the emission altitude: emission altitude and aviation impact. *J Geophys Res Atmos*. <https://doi.org/10.1029/2012JD018204>
- Garfinkel CI, Hartmann DL (2010) Influence of the quasi-biennial oscillation on the North Pacific and El Niño teleconnections. *J Geophys Res Atmos*. <https://doi.org/10.1029/2010JD014181>
- Gottelman A, Hoor P, Pan LL et al (2011) The extratropical upper troposphere and lower stratosphere. *Rev Geophys* 49:RG3003. <https://doi.org/10.1029/2011RG000355>
- Grandpré J, Ménard R, Rochon YJ et al (2009) Radiative impact of ozone on temperature predictability in a coupled chemistry-dynamics data assimilation system. *Mon Wea Rev* 137:679–692. <https://doi.org/10.1175/2008MWR2572.1>
- Grewe V, Dameris M, Fichter C, Lee DS (2002) Impact of aircraft NOx emissions. part 2: effects of lowering the flight altitude. *Meteorol Z* 11:197–205. <https://doi.org/10.1127/0941-2948/2002/0011-0197>
- Haigh JD, Blackburn M, Day R (2005) The response of tropospheric circulation to perturbations in lower-stratospheric temperature. *J Clim* 18:3672–3685. <https://doi.org/10.1175/JCLI3472.1>
- Hall R, Erdélyi R, Hanna E et al (2015) Drivers of North Atlantic polar front jet stream variability. *Int J Climatol* 35:1697–1720. <https://doi.org/10.1002/joc.4121>
- Hegglin MI, Gottelman A, Hoor P et al (2010) Multimodel assessment of the upper troposphere and lower stratosphere: extratropics. *J Geophys Res* 115:D00M09. <https://doi.org/10.1029/2010JD013884>
- Holton JR, Haynes PH, McIntyre ME et al (1995) Stratosphere–troposphere exchange. *Rev Geophys* 33:403–439. <https://doi.org/10.1029/95RG02097>

- Hu S, Vallis GK (2019) Meridional structure and future changes of tropopause height and temperature. *Q J R Meteorol Soc* 145:2698–2717. <https://doi.org/10.1002/qj.3587>
- IPCC AR5, Eds 2014 Climate change 2013—The physical science basis: working group I contribution to the fifth assessment report of the intergovernmental panel on climate change Cambridge University Press Cambridge
- Jacobson MZ, Wilkerson JT, Balasubramanian S et al (2012) The effects of rerouting aircraft around the arctic circle on arctic and global climate. *Clim Change* 115:709–724. <https://doi.org/10.1007/s10584-012-0462-0>
- Jensen EJ, Toon OB, Selkirk HB et al (1996) On the formation and persistence of subvisible cirrus clouds near the tropical tropopause. *J Geophys Res Atmos* 101:21361–21375. <https://doi.org/10.1029/95JD03575>
- Kaminski JW, Neary L, Struzewska J et al (2008) GEM-AQ, an online global multiscale chemical weather modelling system: model description and evaluation of gas phase chemistry processes. *Atmos Chem Phys* 8:3255–3281. <https://doi.org/10.5194/acp-8-3255-2008>
- Kim BY, Fleming GG, Lee JJ et al (2007) System for assessing aviation's global emissions (SAGE), part 1: model description and inventory results. *Transp Res Part D Transp Environ* 12:325–346. <https://doi.org/10.1016/j.trd.2007.03.007>
- Koch P, Wernli H, Davies HC (2006) An event-based jet-stream climatology and typology. *Int J Climatol* 26:283–301. <https://doi.org/10.1002/joc.1255>
- Kuang X, Zhang Y, Huang Y, Huang D (2014) Spatial differences in seasonal variation of the upper-tropospheric jet stream in the Northern Hemisphere and its thermal dynamic mechanism. *Theor Appl Climatol* 117:103–112. <https://doi.org/10.1007/s00704-013-0994-x>
- Lamarque J-F, Emmons LK, Hess PG et al (2012) CAM-chem: description and evaluation of interactive atmospheric chemistry in the community earth system model. *Geosci Model Dev* 5:369–411. <https://doi.org/10.5194/gmd-5-369-2012>
- Lee DS, Fahey DW, Forster PM et al (2009) Aviation and global climate change in the twenty-first century. *Atmos Environ* 43:3520–3537. <https://doi.org/10.1016/j.atmosenv.2009.04.024>
- Lee DS, Pitari G, Grewe V et al (2010) Transport impacts on atmosphere and climate: aviation. *Atmos Environ* 44:4678–4734. <https://doi.org/10.1016/j.atmosenv.2009.06.005>
- Linz M, Chen G, Hu Z (2018) Large-scale atmospheric control on non-gaussian tails of midlatitude temperature distributions. *Geophys Res Lett* 45:9141–9149. <https://doi.org/10.1029/2018GL079324>
- Lund MT, Aamaas B, Berntsen T et al (2017) Emission metrics for quantifying regional climate impacts of aviation. *Earth Syst Dynam* 8:547–563. <https://doi.org/10.5194/esd-8-547-2017>
- Lupu A, Semeniuk K, Kaminski JW, McConnell JC (2013) GEM-AC: a stratospheric-tropospheric global and regional model for air quality and climate change-evaluation of gas-phase properties. In: Bernath PF (ed) *The atmospheric chemistry experiment ACE at 10: a solar occultation anthology*. A. Deepak Publishing, Hampton, Virginia, USA, pp 285–293
- Mamun A, Semeniuk K, Kaminski JW, McConnell JC (2013) Evaluation of stratospheric temperature and water vapor from GEM using ACE-FTS and MLS measurements. In: Bernath PF (ed) *The atmospheric chemistry experiment ACE at 10: A solar occultation anthology*. A. Deepak Publishing, Hampton, Virginia, U.S.A, pp 295–302
- Melamed-Turkish K, Taylor PA, Liu J (2018) Upper-level winds over eastern North America: a regional jet stream climatology. *Int J Climatol* 38:4740–4757. <https://doi.org/10.1002/joc.5693>
- Olsen SC, Brasseur GP, Wuebbles DJ et al (2013a) Comparison of model estimates of the effects of aviation emissions on atmospheric ozone and methane. *Geophys Res Lett*. <https://doi.org/10.1002/2013GL057660>
- Olsen SC, Wuebbles DJ, Owen B (2013b) Comparison of global 3-D aviation emissions datasets. *Atmos Chem Phys* 13:429–441. <https://doi.org/10.5194/acp-13-429-2013>
- Pena-Ortiz C, Gallego D, Ribera P et al (2013) Observed trends in the global jet stream characteristics during the second half of the 20th century. *J Geophys Res Atmos* 118:2702–2713. <https://doi.org/10.1002/jgrd.50305>
- Penner JE (1999) *Aviation and the global atmosphere: a special report of the intergovernmental panel on climate change*. Cambridge University Press
- Rikus L (2018) A simple climatology of westerly jet streams in global reanalysis datasets part 1: mid-latitude upper tropospheric jets. *Clim Dyn* 50:2285–2310. <https://doi.org/10.1007/s00382-015-2560-y>
- Shepherd TG (2007) Transport in the middle atmosphere. *J Meteorol Soc Jpn Ser II* 85B:165–191. <https://doi.org/10.2151/jmsj.85B.165>
- Shepherd TG (2002) Issues in stratosphere-troposphere coupling. *J Meteorol Soc Jpn Ser II* 80:769–792. <https://doi.org/10.2151/jmsj.80.769>
- Simpson IR, Blackburn M, Haigh JD (2012) A mechanism for the effect of tropospheric jet structure on the annular mode-like response to stratospheric forcing. *J Atmos Sci* 69:2152–2170. <https://doi.org/10.1175/JAS-D-11-0188.1>
- Skowron A, Lee DS, De León RR (2015) Variation of radiative forcings and global warming potentials from regional aviation NOx emissions. *Atmos Environ* 104:69–78. <https://doi.org/10.1016/j.atmosenv.2014.12.043>
- Søvde OA, Matthes S, Skowron A et al (2014) Aircraft emission mitigation by changing route altitude: a multi-model estimate of aircraft NOx emission impact on O3 photochemistry. *Atmos Environ* 95:468–479. <https://doi.org/10.1016/j.atmosenv.2014.06.049>
- Strong C, Davis RE (2007) Winter jet stream trends over the Northern Hemisphere. *Q J R Meteorol Soc* 133:2109–2115. <https://doi.org/10.1002/qj.171>
- Sun L, Chen G, Lu J (2013) Sensitivities and mechanisms of the zonal mean atmospheric circulation response to tropical warming. *J Atmos Sci* 70:2487–2504. <https://doi.org/10.1175/JAS-D-12-0298.1>
- Wilkerson JT, Jacobson MZ, Malwitz A et al (2010) Analysis of emission data from global commercial aviation: 2004 and 2006. *Atmos Chem Phys* 10:6391–6408. <https://doi.org/10.5194/acp-10-6391-2010>
- WMO (1957) *Meteorology—a three-dimensional science: Second session of the commission for aerology*. WMO Bull 4:134–138
- Xue D, Lu J, Sun L et al (2017) Local increase of anticyclonic wave activity over Northern Eurasia under amplified arctic warming. *Geophys Res Lett* 44:3299–3308. <https://doi.org/10.1002/2017GL072649>
- Yang H, Waugh DW, Orbe C et al (2019) Large-scale transport into the arctic: the roles of the midlatitude jet and the hadley Cell. *Atmos Chem Phys* 19:5511–5528. <https://doi.org/10.5194/acp-19-5511-2019>
- Zolotov SY, Ippolitov II, Loginov SV (2018) Characteristics of the subtropical jet stream over the North Atlantic from reanalysis data. *IOP Conf Ser Earth Environ Sci* 211:012005. <https://doi.org/10.1088/1755-1315/211/1/012005>





# Relationship between selected percentiles and return periods of extreme events

Dario Camuffo<sup>1</sup> · Francesca Becherini<sup>2</sup> · Antonio della Valle<sup>1</sup>

Received: 20 November 2019 / Accepted: 3 June 2020 / Published online: 10 June 2020  
© Institute of Geophysics, Polish Academy of Sciences & Polish Academy of Sciences 2020

## Abstract

This paper investigates the relationship between selected percentiles, return periods and the concepts of rare and extreme events in climate and hydrological series, considering both regular and irregular datasets, and discusses the IPCC and WMO indications. IPCC (Annex II: Glossary. Contribution of Working Groups I, II and III to the Fifth Assessment Report of the Intergovernmental Panel on Climate Change. IPCC, Geneva, 2014) establishes that an extreme event should be rare and exceed selected upper and lower thresholds (10th and 90th percentiles); WMO (Guidelines on the definition and monitoring of extreme weather and climate events—TT-DEWCE WMO 4/14/2016. World Meteorological Organization, Geneva, 2016) suggests thresholds near the ends of the range, but leaves them undetermined. The concept of “rare” relates the extreme events to the time domain and is typically expressed in terms of return period (RP). The key is to find the combination between “rare”, percentile and return period. In particular, two crucial items are analysed: (1) how the return period may vary in response to the choice of the threshold, in particular when it is expressed in terms of percentiles; (2) how the choice of producing a regular or irregular dataset may affect the yearly frequency and the related return periods. Some weather variables (e.g. temperature) are regular and recorded at fixed time intervals, while other phenomena (e.g. tornadoes) occur at times. Precipitation may be considered either regular, all-days being characterized by a precipitation amount from 0 (no precipitation) to the top of the range, or irregular (rainy-days only) considering a precipitation day over a selected instrumental or percentile threshold. These two modes of interpreting precipitation include a different number of events per year (365 or less) and generate different return periods. Every climatic information may be affected by this definition. The 90th percentile applied to observations with daily frequency produces 10-day return period and the percentiles necessary to get 1 year, 10 years or other return periods are calculated. The general case of events with selected or variable frequencies, and selected percentiles, is also considered with an example of a precipitation series, two-century long.

**Keywords** Extreme events · Return period · Temperature records · Precipitation records · Time series · Percentiles

## Introduction

“Extreme event” is a term nowadays commonly understood and used for a number of different weather phenomena, i.e. heavy precipitation, droughts, earthquakes, tsunamis, epidemics and so forth. Anyhow, its definition is not obvious, or unique, as it is strongly related to the field of application.

In some disciplines, the definition is based only on the magnitude of occurrence of the event, in others, it includes the assessment of its impact on human and natural systems (Broška et al. 2020).

According to WMO (2016), “An extreme can be identified when a single climate variable (e.g. precipitation or wind) exceeds its specific threshold, which can be varying percentile-based values, fixed absolute values and return period”. The definition of extreme events is extremely important, not only for purely academic interest, but also in Earth sciences and in everyday life. It is fundamental not only for research on weather and climate variability (Kharin et al. 2018), but also to assess their adverse effect on landscape and society. For instance, the insurance coverage, the failure of a public service (e.g. transport) or the damage for the collapse of a

✉ Dario Camuffo  
d.camuffo@isac.cnr.it

<sup>1</sup> Institute of Atmospheric Sciences and Climate (ISAC), National Research Council (CNR), Corso Stati Uniti 4, 35127 Padua, Italy

<sup>2</sup> Institute of Polar Sciences (IPS), National Research Council (CNR), Via Torino 155, 3172 Venice Mestre, Italy

structure are reimbursed in a different way (or not at all) if they have been due to a normal or to an exceptional weather situation. In the real world, a number of different approaches have been considered.

In case of stationary systems, it is possible to establish reference levels in absolute terms, i.e. fixed thresholds. In the peak over threshold (POT) theory, all “data exceeding a selected threshold level”, e.g. 30 °C, are considered extreme (Goda 1988, 1992; Allen et al. 2013). This is an intensity-oriented definition, with an absolute threshold.

In the extreme value theory (EVT), extreme events are those “contained in the tail distribution” of a given variable. EVT is tailored on the probability distribution (Galton 1879; Pearson 1895; Fréchet 1927; Gumbel 1941, 1958; Weibull 1951; Tiago de Oliveira 1986; Coles 2001; Katz 2010) of the selected variable and requires an absolute level established as a threshold. The Fréchet distribution is characterized by a heavy-tail with polynomial decay; Gumbel distribution by a double-exponential decay and Weibull distribution by a flexible domain with two parameters. These three approaches can be combined to obtain the generalized extreme value distribution (GEV) (Coles 2001; Neves and Fraga-Alves 2008; Salvadori and De Michele 2013). In EVT, or GEV, the reference definition is: “An extreme (weather or climate) event is generally defined as the occurrence of a value of a weather or climate variable above (or below) a threshold value near the upper (or lower) ends (‘tails’) of the range of observed values of the variable” (WMO 2016).

When climate changes from the condition A to B, also the statistical distribution of the variable changes, and an absolute threshold that was convenient in A will not be longer representative in B. Therefore, the threshold in absolute levels should be updated in terms of moving thresholds, i.e. relative thresholds. A solution is to link the threshold to the distribution, and express it in relative terms, e.g. making reference to a selected percentile level. “An extreme weather event is an event that is rare within its statistical reference distribution at a particular place. Definitions of “rare” vary, but an extreme weather event would normally be as rare as or rarer than the 10th or 90th percentile” (IPCC 2014). This relative threshold will remain unaffected by climate change, even if its absolute value will change.

Percentiles have been popularly adopted for their flexibility and multiple choice that is easily referred to a Gaussian distribution and the standard deviation (SD). For instance, the 50th percentile is the median; the 6.7th and 93.3th correspond to  $\pm 1.5$  SD; the 2.3th and 97.7th to  $\pm 2$  SD; the 0.13th and 99.87th to  $\pm 3$  SD. Alternatively, rounded values of percentiles may be preferred, e.g. 10th and 90th that correspond to  $\pm 1.282$  SD; 1st and 99th to  $\pm 2.326$  SD. The choice of the percentile level is apparently arbitrary; it is clear that the lower or the higher the percentile, the most rare and extreme the event. Several examples are found in

the literature. Osborn et al. (2000) ranked all daily rainfall data, cumulated them and then identified the highest daily amounts that together contribute 10% of the total precipitation, that corresponds to the 90th percentile of the distribution. Miao et al. (2015) used the 95th and the 99th percentiles, as well as absolute threshold indexes (i.e. annual count of days, or the total amount, when precipitation exceeded certain selected thresholds) and Max indices (for some selected annual number of consecutive precipitation) to study extreme precipitation and flood risk in China, with its diverse conditions of geography and topography and its susceptibility to monsoons. The 95th and 99th percentile indices of extreme daily precipitation provided very similar maps over China, but with some differences from the other absolute indexes that were respondent to some specific conditions and different absolute thresholds. Domínguez-Castro et al. (2015) analysed historical extreme precipitation 1855–1856 in Iberia and to this aim they considered various percentile levels (i.e. 10th, 30th, 70th, 90th and 98th) at every Iberian station. The highest percentile levels shown some marked specific features in comparison with the lower ones pointing out some specific periods or areas characterized by abnormal, intense precipitation. Boethe et al. (2018) used selected percentiles to recognize changes in the precipitation distribution in England since 1650 CE. They considered that the 6.7th and 93.3th percentiles (linked to  $\pm 1.5$  SD) are frequently used to represent the regions of severe dryness and wetness, respectively, and decided to compare the evolution of these two percentiles over time, and in addition they considered the 50th percentile as a representative of the average features. A comprehensive review concerning percentiles to assess changes in heavy precipitation can be found in Schär et al. (2016).

The growing interest in trends of extreme weather phenomena is related to their potential for adverse impacts on human life, civil infrastructure and natural ecosystems with socioeconomic consequences (Katz 2010). It is not a case that in some disciplines “extreme events” are also called “natural disasters”, e.g. “Major impacts of climate change on human health are likely to occur via changes in the magnitude and frequency of extreme events which trigger a natural disaster or emergency” (IPCC 2014). The attention to “extremes” is also fostered by concerns that extreme weather and climate events are increasing significantly in frequency and magnitude as a result of global warming while at the same time the natural and human systems are becoming more vulnerable to extremes (Kharin et al. 2018). Although the pathways connecting extreme events to health outcomes and economic losses can be diverse and complex (Stanke et al. 2013), extreme weather and climate-related events affect human health by causing death, physical and mental illness; in addition, they have large socioeconomic impacts (Bell et al 2018). Taking advantage of daily rainfall

measurements in Venezuela, Coles and Pericchi (2003) applied extreme value models to foresee the most critical areas, to be prepared to reduce the impact of disasters with timely interventions and mitigation measures.

Another criterion is based on return period (WMO 2016), also called recurrence interval (AMS 2020). This implies the concept of “exceptional” or “rare” that relates the events to the time domain and is typically expressed in terms of return period (RP) (Elsner and Kara 1999). However, the identification of the concepts of exceptional, extreme and rare should be considered cautiously because they are not always equivalent. For instance, an event may be considered extreme for its impact, but it may not be rare, or vice versa (Yu et al. 2013).

A RP is defined as “a statistical parameter used in frequency analysis as a measure of the average time interval between the occurrence of a given quantity and that of an equal or greater quantity” (Huske 1959) and is expressed as the average time until the next occurrence of a defined event (AMS 2020). The National Academies of Sciences, Engineering and Medicine gives a more detailed definition “a RP is a commonly used metric of probability. If the climate were not changing, RP can also be interpreted as the average time between events, but it should not be interpreted as the time that will pass before an event occurs again” (NASEM 2016). The RP may be interpreted in terms of expectation, and the concept of rarity and RP is particularly useful for insurances, civil engineering and public works, public agencies or rescue teams. Estimating RP implies computing the time intervals between successive events with intensity exceeding a selected threshold (Lestang et al. 2018). Another definition in expectation terms is “in a fixed T-year period, the expected number of exceedances of the T-year event is exactly 1, if the distribution does not change over that period; thus, on average, one event greater than the T-year level occurs in a T-year period” (Stedinger et al 1993). In case of events with successive occurrences independent from one another, the average number of events occurring in a selected time interval is expected to be proportional to the length of that interval and follows a Poisson distribution. When it is necessary to know long RP and related intensities for very rare but catastrophic events, e.g. floods, heavy rainfall, the RP can be estimated using the Poisson distribution of occurrence in samples smaller than the selected RP (Yevjevich and Hatrancioglu 1987). However, it should be considered that the concepts of return periods and return levels are strictly connected to a stationary climate (Katz 2010; Cooley 2013; Salvadori and De Michele 2013) and should be revised in the context of climate change (NASEM 2016; Pendergrass 2018).

In a recent paper (Camuffo et al. 2020a), the most popular definitions of climate and hydrological extremes have been considered and tested on some real case studies. The

relationships between different thresholds, based on standard deviation, percentile, frequency of events exceeding the threshold and return periods were calculated and compared with long temperature and precipitation series. It was found that the 90th percentile threshold suggested by IPCC (2014) gives very short return periods, i.e. 10 days, when applied to long daily temperature or precipitation series, irrespectively of the distribution type and length of the series. Such a short RP was in contrast with the concept of “rare” contained in IPCC (2014) and WMO (2016). Several other examples were presented, but that paper did not include an exhaustive mathematical explanation of how the arbitrary choice of thresholds may affect the returning periods.

This work investigates the relationship between RPs and selected percentiles in climate series as suggested by IPCC (2014). In particular, the relationships between percentiles, return periods, rare events, regular or irregular time series, and length of the dataset are presented. Finally, a theoretical explanation of the findings is provided.

## Difference between regular and irregular time series

The data frequency is connected with percentiles and return periods, but may depend on the definition, that may consider regular or irregular time series, defined as follows: a regular time series stores data for regularly spaced (uniform interval) time points, while an irregular time series stores data for arbitrary time points (nonuniform intervals) (Maidment 2002). The latter should not be confused with series including irregularities, e.g. missing observations, observations collected not regularly over time, or outlying observations (Wright 1986). Some atmospheric variables (e.g. temperature, pressure) are typically regular, flow with continuity and are recorded at fixed time intervals. Other variables have irregular occurrence; some of them are frequent (e.g. precipitation), other are rare (e.g. hailfall, heat waves) or very rare (e.g. floods, tornadoes, volcanic eruptions, earthquakes). For their exceptionality, the wide time window and the irregular interval of years from one event and the subsequent one, the very rare events are generally treated as irregular time series.

The weather phenomena, that occur irregularly, at intervals of days, may be considered in either way. An automatic weather station provides a regular record of all weather variables and their intensities. Concerning precipitation, the observer (or who analyses the data) may decide whether to produce a regular time series composed of all the daily values, which are representative of the observed amount (either zero or different from zero), or to produce an irregular time series composed only of the days in which precipitation had occurred. Basically, one may consider that every time series can be born regular, but

the large number of zeros requires too much memory for storage, and too heavy computation time. The advantage of removing zeros is to reduce the size of the dataset and the calculation power. This is especially relevant in long time series. However, this action transforms a regular series into an irregular one. In other application fields, e.g. computational finance, where market data are typically related to irregular time intervals, it is possible to apply methods (e.g. fast Fourier transform) to resample an irregular time series and transform it into a regular one, although this practice has significant limitations (Song et al. 2014).

Concerning climate or hydrological datasets, the choice regular or irregular series is reversible, because a regular precipitation series may be transformed into an irregular one removing from the dataset all days with zero precipitation; vice versa, an irregular series may be transformed into a regular one, considering all the days of the observing period and attributing amount 0 mm to the days without precipitation. Precipitation is a particularly important variable in climate or hydrological studies and may be interpreted either way, but with different consequences as discussed in the next sections.

## Regular time series

### Frequency of a regular time series

The most typical example of a regular time series is a yearly series composed of 365 daily values and the yearly frequency YF is  $YF = 365 \text{ year}^{-1}$ . This value will be useful in the next formulae, and will be compared with sub- and super-daily frequencies.

Precipitation may be interpreted as a daily or sub-daily variable (Schär et al. 2016), and the series is composed of values that may be greater than zero (rainy days) or zero (no precipitation). Like other regular weather phenomena, all the days of the year, rainy or dry, are represented and contribute to percentiles. The values of zero precipitation are called nil-values (Raes 2013) and do not contribute to the calculation of the monthly or yearly precipitation (either frequency or amounts), but may be used for statistical calculations in climatology or hydrology concerning dry days and aridity periods. Under this point of view, a yearly precipitation record is composed of 365 daily values and  $YF = 365 \text{ year}^{-1}$ , and constitutes a regular series. This approach was used by Moberg et al. (2006), O’Gorman (2014), Ban et al. (2015) and Camuffo et al. (2020a). In the following, other approaches will be presented, which lead to other values of YF. The choice of the approach will be influential on the results, as discussed later.

To make an easier presentation, this and the next sections will proceed by steps, from the simplest case to the most general one.

### The 90th percentile of regular time series of daily values

According to the IPCC (2014) definition, in a series of regular daily variables, in the first year or over the whole period composed of  $n$  years, 90% of the events are “normal” and the complement, i.e. 10%, “extreme”. Therefore, the total number  $N_{\text{tot},l}$  of extreme values is

$$N_{\text{tot},l} = \frac{100 - 90}{100} N_{\text{year}} = 0.1 N_{\text{year}}. \quad (1)$$

The RP of these extreme events is

$$RP = \frac{365}{N_{\text{tot},l}} \text{ (day)}. \quad (2)$$

When 2, 3, ...,  $n$  years are considered, the total time period  $T_{\text{tot}}$  and the total number of extreme events  $N_{\text{tot},n}$  will increase accordingly, i.e.,  $T_{\text{tot}} = n \times 365$  days and  $N_{\text{tot},n} = n \times N_{\text{tot},l}$ , respectively. As a consequence:

$$RP = \frac{n \cdot 365}{N_{\text{tot},n}} = \frac{365}{N_{\text{tot},l}} = 10 \text{ (day)}. \quad (3)$$

The first result is that the RP of extreme events is independent of the length of the observed period. The second is that the events exceeding the 90th percentile are characterized by  $RP = 10$  days. RP is irrespective of the observed record and cannot be used to characterize climate conditions.

Substituting in Eq. (2) the value of  $N_{\text{tot},l}$  given by Eq. (1), Eq. (3) can be rewritten as:

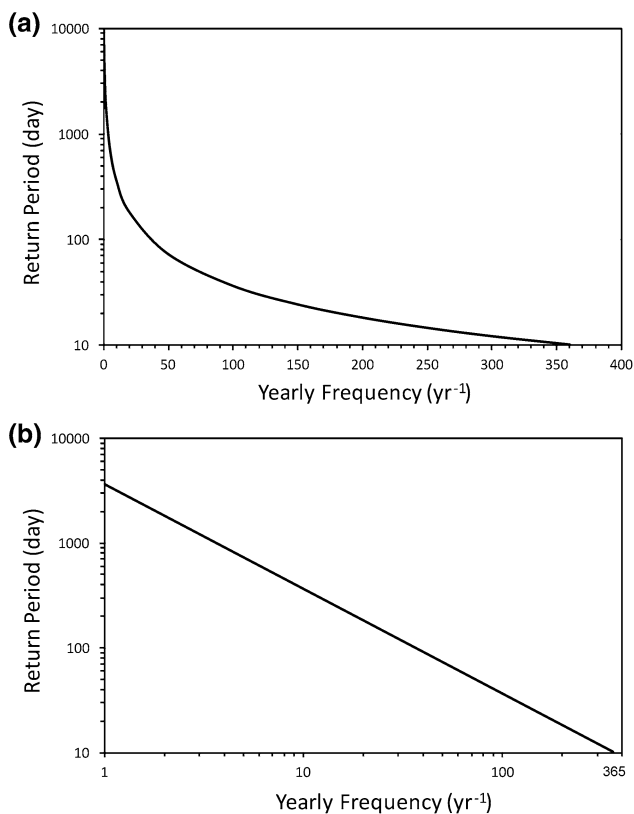
$$RP = \frac{3650}{N_{\text{year}}} = \frac{3650}{YF} = 10 \text{ (day)}, \quad (4)$$

that establishes a direct link between the return period and the yearly frequency of events, i.e. RP is inversely proportional to YF or, which is the same, the variables RP and  $N_{\text{year}}$  are related to each other by a hyperbolic function:

$$RP \times N_{\text{year}} = RP \times YF = 3650 \text{ (day)}. \quad (5)$$

Consequently, in a plot of Eqs. (4) and (5) (Fig. 1a), the abscissa may be exchanged with the ordinate.

Long daily temperature series constitute a typical example of the case discussed in this section. Equation (4) gives  $RP = 10$  days and the same result was empirically obtained in Camuffo et al. (2020a) by counting the number of events that exceeded the 90th percentile in a long series of daily temperatures, i.e. Bologna from 1715 to 2016 (Camuffo et al. 2017). The same result was obtained with the long series of



**Fig. 1** How the 90th percentile matches the return period with the number of daily events per year. **a** Log-linear plot; **b** log–log plot. Calculation made for regular time series with daily values

daily (regular type) precipitation, i.e. Bologna from 1813 to 2016 (Brunetti et al. 2001; Camuffo et al. 2019).

For every selected year, values exceeding the 90th percentile obtained with the above definition of RP (Huske 1959) are trivial and do not constitute climate information. Every year will have ten extreme hot days in summer (as well as 10 extreme cold days in winter). However, the situation is different if the selected percentile is calculated over a certain number of years. A time series composed of  $n$  years will have  $10 \times n$  extreme hot days (as well as  $10 \times n$  extreme cold days), but these will likely differ from the highest 10 of each year (as well from the 10 lowest of each year). Therefore, the set of days exceeding the selected percentile constitutes a new series of irregularly distributed events (Lestang et al., 2018). These “extreme” days may be evenly distributed over the whole period (every year with 10 hot days = stationary climate) or may be concentrated in some warmer years or warmer periods; the years (or periods) with less of these extreme days being considered colder. The climate information is how the extremes are distributed over time (e.g. trends) and space (e.g. characterization of regional climate). As an example, Przybylak et al. (2007) considered for Poland the trends of some climate variables, including the

yearly frequency of days in which the maximum temperature exceeded the 90th and 99th percentiles calculated over the 1950–2005 period. The result was that the frequency of the hot days exceeding the 90th percentile level increased with slope 0.36 hot days/year, and at the 99th level with 0.09 hot days/year. The indication was global warming, with larger increase of days exceeding the first threshold, and lower exceeding the second threshold. This example shows two (obvious) conclusions. The result depends on (1) the choice of the thresholds; (2) the peculiarity of the dataset. The analysis has extracted the information included in the dataset: the trend was evident as the dataset was long enough to include the recent climate change. In fact, if the period had been 30 years, from 1990 to 2020, a more stationary situation would have appeared. If the analysis had been made considering different, indirect parameters, e.g. crop production or socioeconomic impacts, probably the result could have been different, because for their complexity, not all variables, and not all extremes can be approached with the same statistical analysis (Coles 2001; Neves and Fraga-Alves 2008).

### The 90th percentile of sub- and super-daily frequency of regular time series

The general case with a number of values smaller or greater than 1 per day is similar to the previous section. The formulae are the same, but the number of events per year should be changed according to their frequency, e.g.

2 values per day,  $N_{\text{year}} = YF = 2 \times 365$ , hence  $RP = 5$  days;

1 value every 2 days,  $N_{\text{year}} = YF = 365/2$ , hence  $RP = 20$  days;

100 values per year,  $N_{\text{year}} = YF = 100$ , hence  $RP = 36.5$  days.

Plots of matched values of YF and RP are reported in Fig. 1a, b.

### Percentiles in a regular time series of daily values

In case of daily observations, i.e.  $N_{\text{year}} = 365 \text{ year}^{-1}$ , and percentiles higher than a selected value, Eq. (1) becomes

$$N_{\text{tot},1} = \frac{100 - SP}{100} N_{\text{year}} = PT \times N_{\text{year}}, \quad (6)$$

where SP is the selected percentile, e.g.  $SP = 95$  for the 95th percentile. To this aim, the selected normalized percentile threshold PT may be defined as

$$PT = \frac{100 - SP}{100}. \quad (7)$$

The key Eq. (5) assumes the general form:

$$RP \times N_{\text{year}} = \frac{3650}{100 - SP} = \frac{365}{PT}. \quad (8)$$

Equation (2) can be rewritten

$$RP = \frac{365 \times 100}{N_{\text{year}}(100 - SP)} = \frac{100}{100 - SP} = \frac{1}{PT} \quad (\text{day}). \quad (9)$$

With this equation, it is possible to calculate specific values of SP (Fig. 2a, b) to get particular return periods.

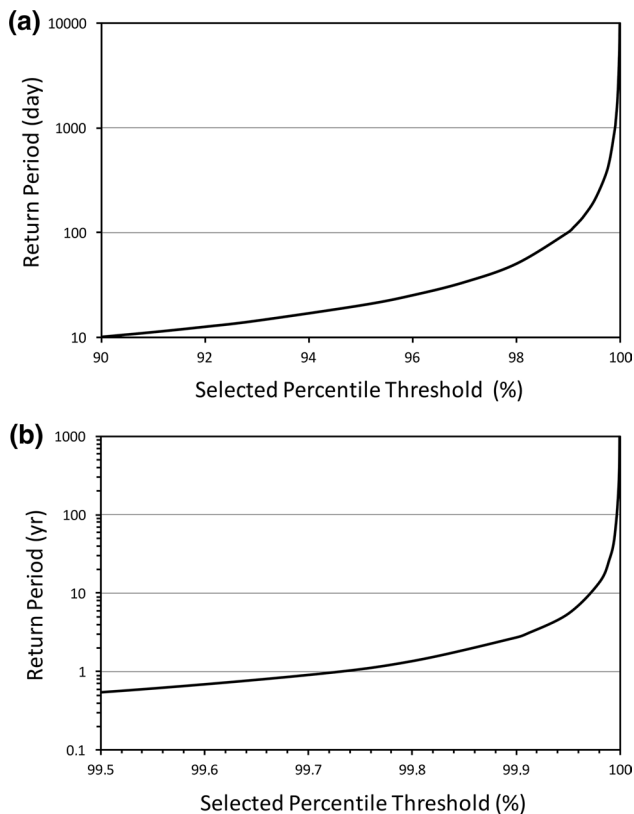
For instance,  $RP = 1$  year is obtained when the denominator of Eq. (9) is  $100 - SP = 100/365$  and this particular percentile  $SP_{1\text{yr}}$  is

$$SP_{1\text{yr}} = 99.7260274 \quad (\text{percentile}). \quad (10)$$

To get  $RP > 1$  year, it is necessary to consider higher percentiles.

Similarly, for events with  $RP = 10$  year, the needed percentile  $SP_{10\text{yr}}$  is

$$SP_{10\text{yr}} = 99.9726028 \quad (\text{percentile}). \quad (11)$$



**Fig. 2** Return periods for selected percentile thresholds from 90 to 100th percentile. **a** RP in day for SP in the 90–100th percentile range. **b** RP in years for SP in the 99.5–100th percentile range. Calculation made for regular time series with daily values

## Irregular time series

### Frequency of irregular time series

In irregular time series, the data sequence is not regular and every year is composed by a variable number of values, from a few to several ones. Precipitation may be interpreted as an event-oriented variable, i.e. only rainy days contribute to percentiles, while days without precipitation are excluded from the series. Precipitation is considered only when the rain gauge gives an output different from zero, i.e. the precipitated water exceeds the instrumental threshold, e.g., 0.1 mm (Sneeyers 1990; Rajczak et al. 2013; Kendon et al. 2014). Another similar criterion is based on a percentile threshold (Zhang et al. 2011; Sillmann et al. 2013; Giorgi et al. 2014).

The series used as examples will be characterized by a number  $N_{\text{year}}$  of events per year and  $N_{\text{tot},I}$  will be the total number of “extreme values” over a selected period, i.e. the values above a selected threshold. For every individual year  $i$ , their frequency  $YF_i$  is  $YF_i < 365 \text{ year}^{-1}$ , with  $YF_i$  being different year by year. In the case of a long series composed of  $n$  years, one should consider the average yearly frequency  $YF$  over the whole series, i.e.,

$$YF = \sum_i \frac{YF_i}{n}. \quad (12)$$

The instrumental or percentile threshold is a crucial issue, because in early series, it is often unknown and cuts off the lowest amounts up to a certain unknown limit. As the probability density function of daily precipitation amounts follows a Gamma, Weibull or double-exponential function or other functions (Wilks 2011, Schär et al. 2016), in a precipitation series, the lowest amounts are the most frequent ones. The combination of records taken with different instruments (i.e. with different thresholds) may constitute a serious bias (Camuffo et al. 2019, 2020b). The uncontrolled cut-off of the lowest amount is especially relevant for frequency, less for the total amount. If the interest is focused on the most intense precipitation, it is recommended to cut-off the low precipitation readings that are the ones more affected by the instrumental threshold.

### The 90th percentile of irregular time series

In irregular time series composed of event-oriented variables, the dataset may be one or several years long. The treatment for frequencies lower than 365 per year is similar to the case of the regular daily frequencies, and is based on the equation:

$$RP = \frac{3650}{YF} \text{ (day)}, \quad (13)$$

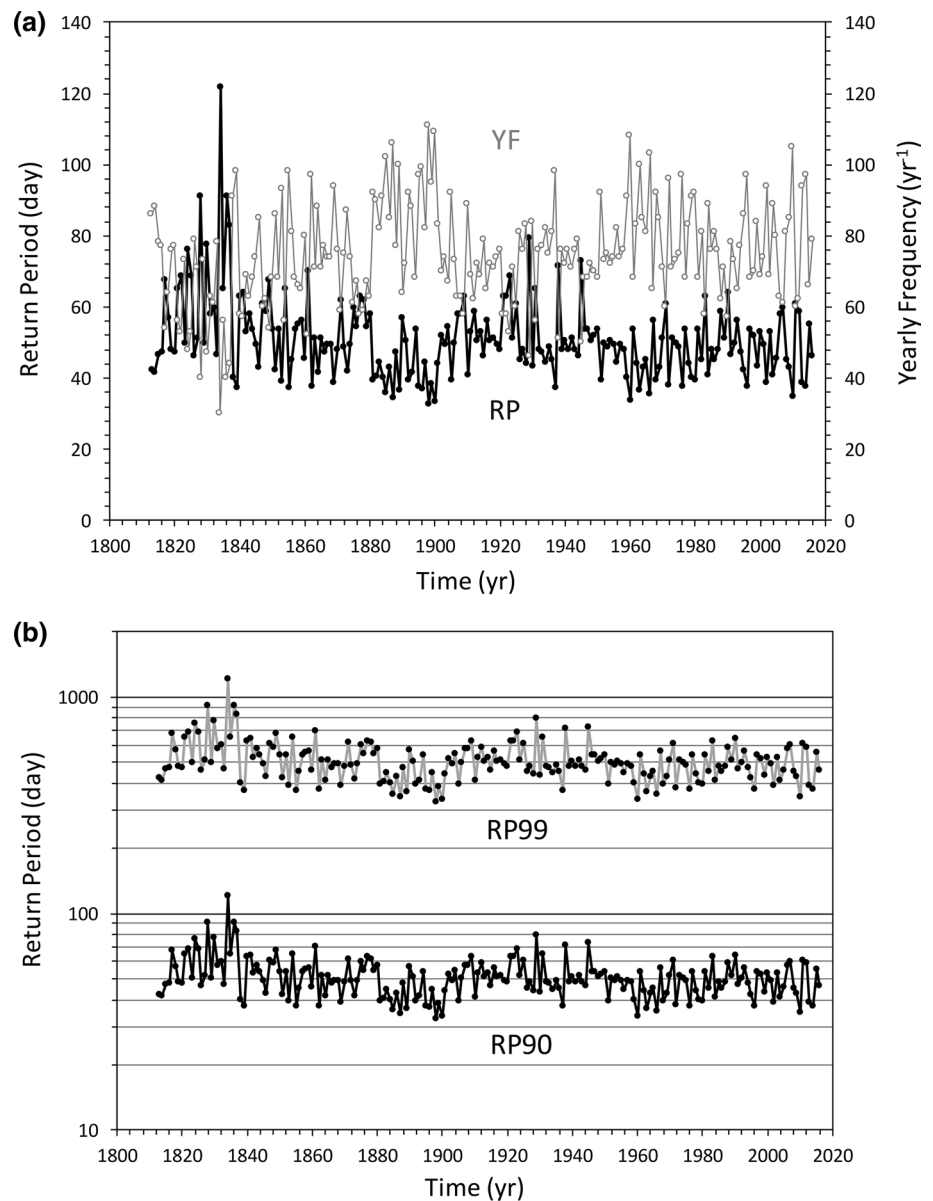
where YF is now variable and related to climate. It represents the average yearly frequency when more than 1 year is considered. From this formula, it is evident that the lower YF, the longer RP, and both variables are equally useful to point out climate changes or particular climate periods.

An example of only rainy-day precipitation type with threshold 1 mm is given in Fig. 3a, b, where the series of daily precipitation in Bologna, Italy, from 1813 to 2016, has been used (Brunetti et al. 2001). The plot provides information about the climate variability, e.g. periods in which the increased precipitation frequency is associated with shorter return periods and vice versa. Apparently, the lowest RP,

or the highest YF, were in the periods 1880–1900 and then 1960–1980, while the highest RP, or the lowest YF, in 1820–1840 and 1900–1940.

However, as for every long series, the plot may be affected by instrumental bias when an instrument has been changed, because every instrument is characterized by a threshold that cuts off the lighter precipitation and the probability density function of precipitation reaches the highest values at the lowest amounts. The resulting precipitation distribution is not affected at intermediate to high percentiles, but may be significantly affected at low percentiles and in terms of rainy-day frequency. This is a serious problem in the long series, because the cut-off value of early instruments is generally unknown and this instrumental bias risks to be misinterpreted for climate signal (Camuffo et al. 2019, 2020b)

**Fig. 3** **a** Return period of precipitation events above the 90th percentile (RP, black) and yearly frequency (YF, grey) of the Bologna series from 1813 to 2016. Irregular time series composed of daily amounts higher than 1 mm. **b** Return period of events above the 90th and 99th percentile (RP90 and RP99, respectively) for the above series



even if various techniques may be applied to detect discontinuities and homogenize series affected by instrumental, location or other changes (Craddock 1979; Peterson et al. 1998; Aguilar et al. 2003; Wijngaard et al. 2003; Costa and Soares 2009; Todd et al. 2015; WMO 2018). The climate signal becomes more evident if the analysis is performed considering the medium and intense precipitation and disregarding the light one, e.g. < 1 mm per day, as in Fig. 3a, b. This confirms the observation that rainy-day percentiles are very sensitive to thresholds and evaporation losses that may affect the fraction of rainy days as well as their distribution, and may produce misleading results when used to address changes in heavy precipitation events (Schär et al. 2016).

### Percentiles related to an irregular time series

The most general case is with any yearly frequency YF and any selected percentile threshold PT. The basic equation becomes

$$RP = \frac{365}{YF \times PT} = \frac{36500}{YF(100 - SP)} \quad (\text{day}), \quad (14)$$

where the variables are two: YF related to climate; PT to an arbitrary choice of the percentile threshold. This equation is very general and may be applied to the regular events too, using the appropriate YF. From this equation, it is evident that the return periods of SP having the particular values 90; 99; 99.9; 99.99; 99.999 etc. are related between them as 1; 10; 10<sup>2</sup>; 10<sup>3</sup>; 10<sup>4</sup> and so forth.

An example for SP = 90 and 99 is shown in Fig. 3b for Bologna. Passing from SP = 90 to 99, RP increases by an order of magnitude. The two RP trends, however, remain unchanged. In this particular series and with 1 mm threshold, a certain variability was visible especially in the first half of the nineteenth century, when the rain gauge was located on the Astronomical Tower at 48 m above ground level, thus strongly influenced by the more intense and variable wind field (Respighi, 1857; Camuffo et al. 2019). Milestones in instrument and exposure standardization were in 1865, when the first Italian weather service for navy (named *Servizio Operativo Marittimo*) was created; 1879 when a general weather service (named *Ufficio Centrale di Meteorologia*) supervised all stations; 1925 when the forecast weather service of the air force (named *Ufficio Presagi*) became operative. Over the whole twentieth century, when the rain gauge was at ground level under standard conditions, this parameter became stable, without evidence of increased frequency of extreme events over the most recent decades. This shows how the data homogeneity is a crucial prerequisite in long series.

Finally, this explains why Gilleland and Katz (2011) found the 95th and especially the 99th percentiles strongly

correlated with the decadal top ten precipitation totals; similarly, Tu and Chou (2013), Knapp et al. (2015), Salack et al. (2018) and Wasko et al. (2018) preferred to classify as “extreme events” the most intense events on the basis of the 99th percentile.

### Conclusions

This paper has analysed two crucial items: (1) how the return period may vary in response to the arbitrary choice of the threshold and, in particular, when it is expressed in terms of percentiles; (2) how the choice of producing a dataset in form of regular, or irregular time series may affect the yearly frequency and the related return periods.

The RP of an event is related to the occurrence, or the observation frequency, of a selected weather phenomenon. Depending on the weather variable, time series may be composed of values distributed at regular time intervals (continuous variables like temperature and pressure), or at irregular ones (event-oriented variables like extreme or less common weather events like floods, tornadoes and so on).

Precipitation constitutes a particular variable because it may be considered from two different points of view: (1) regular daily precipitation amounts, ranging from 0 (no precipitation) to the highest amount of collected water; (2) irregular daily precipitation amounts, starting from a selected instrumental or percentile threshold (i.e. only rainy days). The latter has lower YF and gives longer RP. Any quantitative information about climate changes is affected by the definition chosen.

When regular time series are considered, RP is irrespective of the observed record and is uniquely determined by the selected percentile threshold. RP and YF are inversely proportional between them and the RP of events exceeding selected percentiles is independent of the length of the series. As RP depends on the arbitrary choice of the percentile, its value is determined by the percentile and is not useful to characterize a particular climate. However, different climate periods will be characterized by an uneven distribution of such extremes and this distribution may constitute climate information.

When irregular time series are considered (e.g. rainy-days only), RP depends on the selected percentile as well as on the actual value of YF, and may characterize a particular climate period.

In the general case of events at sub- or super-daily time scale and higher percentiles, RP is found considering that it is inversely proportional to the product of the yearly frequency by the selected percentile level.

In the real world, WMO (2016) suggests a threshold value closer to the ends of the range; however, the value is not specified. IPCC (2014) suggests a lower and an upper



threshold (i.e. 10th and 90th percentile) that may be convenient for some datasets and certain purposes, but not for all. The results of this paper are in accordance with other Authors, i.e.: the approach is not unique, especially when dealing with variables that may be defined in different ways, e.g. extreme precipitation (Pendergrass 2018). In addition, the choice of the statistical approach may be influential on the result (e.g. Coles 2001; Neves and Fraga-Alves 2008).

This paper gives the key to calculate a proper percentile threshold for any observational frequency and any selected RP, and has shown that, in the case of regular series, the 90th percentile may give too short RPs, which dissociates the concept of extreme from the concept of rare. In particular, it has been found that for a daily series it is necessary to pass from the 90th to the 99.9726028th percentile to move from a 10-day to a 10-year RP. This poses the question: what is the most convenient length of a RP for an event to be considered “rare” and therefore extreme, as suggested by WMO (2016)? In the real world, the concept of rare should be tailored on the subject that should be protected from extreme events. For example, the agriculture is vulnerable to moderately severe events and focuses on short RPs, while the infrastructure sector focuses on very long RPs, because civil works like bridges and dams should resist to extreme events with very low probability of occurrence.

This study suggests in particular that percentile thresholds and RPs are strictly related between them by mathematical formulae, but their concepts are not equally related, because relatively high percentiles may be associated with relatively short RPs. In other words, it may be misleading to use them synonymously, i.e. high percentile equal to long RP, but each of them has an individual meaning and should be calculated for specific aims correspondent to their very definitions. Briefly, in some cases, the percentile definition may be preferable (e.g. the analysis of a meteorological or hydrological records), while in others, the RP approach (e.g. forecast or design of engineering structures).

**Acknowledgements** The authors are grateful to the two anonymous Referees for the useful suggestions; to Michele Brunetti (CNR-ISAC) and i.e. the Agenzia Regionale per la Prevenzione e Protezione Ambientale dell’Emilia-Romagna (ARPA ER), for having kindly provided precipitation data.

**Funding** This study did not receive any specific grant from funding agencies in the public, commercial or not for profit sectors.

### Compliance with ethical standards

**Conflict of interest** The authors declare that they have no conflict of interest.

## References

- Aguilar E, Auer I, Brunet M, Peterson TC, Wieringa J (2003) Guidelines on climate metadata and homogenization. World Meteorological Organization, WMO-TD No. 1186, WCDMP No. 53, Geneva
- Allen DE, Singh AK, Powell RJ (2013) EVT and tail-risk modelling: evidence from market indices and volatility series. *North Am J Econ Finan* 26:355–369
- AMS (2020) Meteorology glossary. American Meteorological Society, Boston. [https://glossary.ametsoc.org/wiki/Main\\_Page](https://glossary.ametsoc.org/wiki/Main_Page)
- Ban N, Schmidli J, Schär C (2015) Heavy precipitation in a changing climate: does short-term summer precipitation increase faster? *Geophys Res Lett* 42:1165–1172. <https://doi.org/10.1002/2014GL062588>
- Bell JE, Brown CL, Conlon K, Herring S, Kunkel KE, Lawrimore J, Luber G, Schreck C, Smith A, Uejji C (2018) Changes in extreme events and the potential impacts on human health. *J Air Waste Manage* 68(4):265–287
- Bothe O, Wagner S, Zorita E (2018) Inconsistencies between observed, reconstructed, and simulated precipitation indices for England since the year 1650 CE. *Clim Past* 15:307–334
- Broska LH, Poganietz WR, Vogeles S (2020) Extreme events defined—a conceptual discussion applying a complex systems approach. *Futures* 115:102490
- Brunetti M, Buffoni L, Lo Vecchio G, Maugeri M, Nanni T (2001) Tre secoli di meteorologia a Bologna. CUSL, Milan
- Camuffo D, della Valle A, Bertolin C, Santorelli E (2017) Temperature observations in Bologna, Italy, from 1715 to 1815; a comparison with other contemporary series and an overview of three centuries of changing climate. *Clim Chang* 142(1–2):7–22
- Camuffo D, Becherini F, della Valle A (2019) The Beccari Series of Precipitation in Bologna, Italy, from 1723 to 1765. *Clim Chang* 155:359–376
- Camuffo D, della Valle A, Becherini F (2020a) A critical analysis of the definitions of climate and hydrological extreme events. *Quat Int* 538:5–13
- Camuffo D, Becherini F, della Valle A (2020b) Three centuries of daily precipitation in Padua, Italy, 1713–2018. History, relocations, gaps, homogeneity and raw data. *Clim Chang*. <https://doi.org/10.1007/s10584-020-02717-2>
- Coles S (2001) An introduction to statistical modeling of extreme values. Springer Series in Statistics, London
- Coles S, Pericchi L (2003) Anticipating catastrophes through extreme value modelling. *J R Stat Soc Ser C* 52(4):405–416
- Cooley D (2013) Return periods and return levels under climate change. In: AghaKouchak A, Easterling D, Hsu K, Schubert S, Sorooshian S (eds) *Extremes in a changing climate: detection, analysis and uncertainty*. Springer, Dordrecht, pp 97–114
- Costa AC, Soares A (2009) Homogenization of climate data: review and new perspectives using geostatistics. *Math Geosci* 41:291–305
- Craddock JM (1979) Methods of comparing annual rainfall records for climatic purposes. *Weather* 34(9):332–346
- Domínguez-Castro F, Ramos AM, García-Herrera R, Trigo RM (2015) Iberian extreme precipitation 1855/1856: an analysis from early instrumental observations and documentary sources. *Int J Climatol* 35:142–153
- Du T, Xiong L, Xu CY, Gippel CJ, Guo S, Liu P (2015) Return period and risk analysis of nonstationary low-flow series under climate change. *J Hydrol* 527:234–250
- Elsner JB, Birol Kara A (1999) *Hurricanes of the North Atlantic: climate and society*. Oxford University Press, New York
- Fréchet M (1927) Sur la loi de probabilité de l’écart maximum. *Ann Soc Polon Math* 6:93–116

- Galton F (1879) The geometric mean in vital and social statistics. *Proc R Soc Lond* 29:365–367
- Gilleland E, Katz RW (2011) New software to analyze how extremes change over time. *Eos Trans AGU* 92(2):13–14
- Giorgi F, Coppola E, Raffaele F (2014) A consistent picture of the hydroclimatic response to global warming from multiple indices: models and observations. *J Geophys Res Atmos* 119:11695–11708
- Goda Y (1988) On the methodology of selecting design wave height. In: *Proc. 21st International Conference on Coastal Engineering*. Malaga, pp 899–913
- Goda Y (1992) Uncertainty of design parameters from viewpoint of extreme statistics. *J Offshore Mech Arct Eng* 114(2):76–82
- Gumbel EJ (1941) The return period of flood flows. *Ann Math Statist* 12(2):163–190
- Gumbel EJ (1958) *Statistics of Extremes*. Columbia University Press, New York
- Huske RE (1959) *Glossary of meteorology*. American Meteorological Society, Boston
- IPCC (2014) *Climate change 2014: synthesis report*. In: Mach KJ, Meyer LA, Pachauri RK, Planton S, von Stechow C (eds) Annex II: Glossary. Contribution of Working Groups I, II and III to the Fifth Assessment Report of the Intergovernmental Panel on Climate Change. IPCC, Geneva
- Katz RW (2010) Economic impact of extreme events: an approach based on extreme value theory. In: Chavez M, Ghil M, Urrutia-Fucugauchi J (eds) *Extreme events: observations, modeling, and economics*. AGU Wiley, Hoboken, pp 207–217
- Kendon EJ, Roberts NM, Fowler HJ, Roberts MJ, Chan SC, Senior CA (2014) Heavier summer downpours with climate change revealed by weather forecast resolution model. *Nat Clim Chang* 4:570–576
- Kharin VV, Flato G, Zhang X, Gillett NP, Zwiers F, Anderson KJ (2018) Risks from climate extremes change differently from 1.5 °C to 2.0 °C depending on rarity. *Earths Future* 6(5):704–715
- Knapp AK, Hoover DL, Wilcox KR, Avolio ML, Koerner SE, Kimberley J, La Pierre KJ, Loik ME, Luo Y, Sala OE, Smith MD (2015) Characterizing differences in precipitation regimes of extreme wet and dry years: implications for climate change experiments. *Global Change Biol* 21:2624–2633
- Lestang T, Ragone F, Bréhier CE, Herbert C, Bouchet F (2018) Computing return times or return periods with rare event algorithms. *J Stat Mech* 2018:043213
- Maidment DR (2002) Arc hydro GIS for water resources, chapter 7 time series. ESRI Press, Redland
- Miao C, Ashouri H, Hsu KL, Sorooshian S, Duan Q (2015) Evaluation of the PERSIANN-CDR daily rainfall estimates in capturing the behavior of extreme precipitation events over China. *J Hydrometeorol* 16:1387–1396
- Moberg A, Jones PD, Lister D, Walther A, Brunet M, Jacobeit J, Alexander LV, Della-Marta PM, Luterbacher J, Yiou P, Chen D, Klein Tank AMG, Saladie O, Sigro J, Aguilar E, Alexandersson H, Almarza C, Auer I, Barriendos M, Begert M, Bergström H, Böhm R, Butler CJ, Caesar J, Drebs A, Founda D, Gerstengarbe FW, Micela G, Maugeri M, Osterle H, Pandzic K, Petrakis M, Srnc L, Tolasz R, Tuomenvirta H, Werner PC, Linderholm H, Philipp A, Wanner H, Xoplaki E (2006) Indices for daily temperature and precipitation extremes in Europe analyzed for the period 1901–2000. *J Geophys Res* 111(D22106):1–25
- NASEM (National Academies of Sciences, Engineering, and Medicine) (2016) *Attribution of Extreme Weather Events in the Context of Climate Change*. The National Academic Press, Washington, DC
- Neves C, Fraga-Alves MI (2008) Testing extreme value conditions—an overview and recent approaches. *Revstat* 6:83–100
- O’Gorman PA (2014) Contrasting responses of mean and extreme snowfall to climate change. *Nature* 512:416–418
- Osborn TJ, Hulme M, Jones PD, Basnett TA (2000) Observed trends in the daily intensity of United Kingdom precipitation. *Int J Climatol* 20:347–364
- Pearson K (1895) Contributions to the mathematical theory of evolution. II. Skew variation in homogeneous material. *Philos Trans R Soc Lond A* 186:343–414
- Pendergrass A (2018) What precipitation is extreme? *Science* 360(6393):1072–1073
- Peterson TC, Easterling DR, Karl TR, Groisman P, Nicholls N, Plummer N, Torok S, Auer I, Böhm R, Gullett D, Vincent L, Heino R, Tuomenvirta H, Mestre O, Szentimrey T, Salinger J, Førland EJ, Hanssen-Bauer I, Alexandersson H, Jones PE, Parker D (1998) Homogeneity adjustments of in situ atmospheric climate data: a review. *Int J Climatol* 18:1493–1517
- Przybylak R, Vizi Z, Arażny A, Kejna M, Maszewski R, Uscka-Kowalkowska J (2007) Poland’s Climate Extremes Index, 1951–2005. *Geogr Pol* 80:47–58
- Raes D (2013) Frequency analysis of rainfall data. College on Soil Physics- 30th Anniversary (1983–2013). The Abdus Salam International Centre for Theoretical Physics, Trieste, pp 1–42
- Rajczak J, Pall P, Schär C (2013) Projections of extreme precipitation events in regional climate simulations for Europe and the Alpine region. *J Geophys Res Atmos* 118:3610–3626
- Respighi L (1857) Notizie sul clima bolognese. *Rendiconto delle sessioni della R. Accademia delle Scienze dell’Istituto di Bologna*. San Tomaso d’Aquino, Bologna
- Salack S, Saley IA, Bliefernicht J (2018) Observed data of extreme rainfall events over the West African Sahel. *Data Brief* 20:1274–1278
- Salvadori G, De Michele C (2013) Multivariate extreme value methods. In: Aghakouchak A, Easterling D, Hsu K, Schubert S, Sorooshian S (eds) *Extremes in a changing climate: detection, analysis and uncertainty*. Springer, Dordrecht, pp 115–162
- Schär C, Ban N, Fischer EM, Rajczak J, Schmidli J, Frei C, Giorgi F, Karl TR, Kendon EJ, Klein Tank AMG, O’Gorman PA, Sillmann J, Zhang X, Zwiers FV (2016) Percentile indices for assessing changes in heavy precipitation events. *Clim Chang* 137:201–216
- Sillmann J, Kharin VV, Zwiers FW, Zhang X, Bronaugh D (2013) Climate extremes indices in the CMIP5 multimodel ensemble: part 2. Future climate projections. *J Geophys Res Atmos* 118:1–21
- Sneyers R (1990) On the statistical analysis of series of observations. Technical note WMO No 415. World Meteorological Organization, Geneva
- Song JH, de Prado ML, Simon HD, Wu K (2014) Exploring irregular time series through non-uniform fast Fourier transform. In: *WHPCF ’14: Proceedings of the 7th Workshop on High Performance Computational Finance*, IEE Computer Society, pp 37–44
- Stanke C, Kerac M, Prudhomme C, Medlock J, Murray V (2013) Health effects of drought: a systematic review of the evidence. *PLoS Curr*. <https://doi.org/10.1371/currents.dis.7a2cee9e980f91ad7697b570bcc4b004>
- Stedinger J, Vogel R, Foufoula-Georgiou E (1993) Frequency analysis of extreme events. In: Maidment D (ed) *Handbook of hydrology*. McGraw-Hill, New York, pp 181–1868
- Tiago de Oliveira J (1986) Extreme values and meteorology. *Theor Appl Climatol* 37(4):184–193
- Todd B, Macdonald N, Chiverrell RC (2015) Revision and extension of the composite Carlisle rainfall record, northwest England: 1757–2012. *Int J Climatol* 35:3593–3607
- Tu JY, Chou C (2013) Changes in precipitation frequency and intensity in the vicinity of Taiwan: typhoon versus non-typhoon events. *Environ Res Lett* 8:014023
- Wasko C, Lu WT, Mehrotra R (2018) Relationship of extreme precipitation, dry-bulb temperature, and dew point temperature across Australia. *Environ Res Lett* 13:074031
- Weibull W (1951) A statistical distribution function of wide applicability. *J Appl Mech Trans ASME* 18(3):293–297

- Wijngaard J, Klein Tank AMG, Können GP (2003) Homogeneity of 20th century European daily temperature and precipitation series. *Int J Climatol* 23(6):679–692
- Wilks DS (2011) *Statistical methods in the atmospheric sciences*. International geophysics series, 3rd edn. Academic Press, Oxford
- WMO (2016) *Guidelines on the Definition and Monitoring of Extreme Weather and Climate Events—TT-DEWCE WMO 4/14/2016*. World Meteorological Organization, Geneva
- WMO (2018) *Guidance on the homogenization of climate station data*. EarthArXiv World Meteorological Organization, Geneva
- Wright DJ (1986) Forecasting data published at irregular time intervals using extension of Holt's method. *Manage Sci* 32(4):499–510
- Yevjevich V, Hatrmancioglu NB (1987) Research needs on flow characteristics. In: Singh VP (ed) *Application of frequency and risk in water resources*. Reidel, Dordrecht, pp 1–22
- Yu T, Chawla N, Simoff S (2013) *computational intelligent data analysis for sustainable development*. CRC Press, Boca Raton
- Zhang X, Alexander L, Hegerl GC, Jones PD, Klein Tank A, Peterson TC, Trewin B, Zwiers FW (2011) Indices for monitoring changes in extremes based on daily temperature and precipitation data. *WIREs Clim Chang* 2:851–870



# Site diversity performance in Ka band using a 7.3-m antenna diameter at tropical climate: a comparison of prediction models

Fazdliana Samat<sup>1</sup> · Mandeep Jit Singh<sup>1,2</sup>

Received: 25 June 2019 / Accepted: 3 June 2020 / Published online: 9 June 2020  
© Institute of Geophysics, Polish Academy of Sciences & Polish Academy of Sciences 2020

## Abstract

Site diversity gain prediction models were created to estimate mathematically the acquired benefits from the implementation of site diversity at place of choice. This work contributes to the comparison of existing gain prediction model to the gain of measured attenuation at Cyberjaya and Rawang, Malaysia. The experiment has been conducted for 4 years from 2014 to 2017, in Ka band using a large 7.3-m diameter antenna and a high elevation angle of 68.8°, together with the rain analysis at both places for the same duration. The average monthly rainfall and attenuation for 4 years were presented. The results revealed that prediction model Hodge performs better than other models, while X. Yeo and Panagopoulos models appear to exhibit very similar graph shape to the measured gain data. More research on gain development in tropical region should be conducted, as the existing prediction model appears to be less consistent with the current data.

**Keywords** Signal propagation · Atmospheric attenuation · Ka-band signal · Model comparison · Site diversity

## Introduction

The trend of wireless communication is gearing to 5G network, which has capabilities of higher speed and larger bandwidth than earlier technology (4G). While terrestrial network is becoming a focus, satellite communication plays similar significant role to provide the requested future bandwidth capabilities, especially at unreachable area by the former. Ka band, Q/V band or even W band is a band of choice to provide this demand, because higher frequency correlates with higher capacity (Kyrgiazos et al. 2014). With high throughput satellite (HTS), this high capacity could be realized and resulted in consequent cost saving at roughly 30% per annum in view of satellite communication service provider (Callaghan et al. 2008). Unfortunately, the high

frequency is easily being degraded by tropospheric precipitation and scintillation. However, the most impairment comes from rain, together with cloud formation that contributed to the attenuation of the signal. Nonetheless, the effect of cloud attenuation is seen around 2–4 dB, impacting the received signal from the satellite (Omotsho et al. 2011; Yuan et al. 2017). This formation of cloud and the correlated attenuation is observed severe for VSAT (very small aperture terminal) (Yuan et al. 2016) at longer slant path and low elevation angle (Yang et al. 2013). For the past years, many techniques have been proposed to mitigate this weather effects. Most of the experiments studied the effect of rain attenuation at Ka-band frequency, using temperate region databank. The experimental samples observed in this region showed that the attenuation at Ka band is more severe than Ku, at mostly double the effect (Panagopoulos et al. 2004).

Fade mitigation technique (FMT) proposed power control, adaptive coding modulation (ACM) and diversity techniques, as solutions depending on the purpose of the system application and requirement (Yussuff et al. 2017). However, since higher frequencies are expected to experience higher degree of signal degradations, so power control and ACM could no longer be applied (Rytir et al. 2017). Diversity itself is divided into four ways, namely frequency, time, satellite and site diversity. Among these techniques, site diversity (SD) is viewed as more efficient (Ippolito 2017). The

✉ Fazdliana Samat  
fazdliana@yahoo.com

Mandeep Jit Singh  
mandeep@ukm.edu.my

<sup>1</sup> Department of Electrical, Electronic and System Engineering, Faculty of Engineering and Built Environment, Universiti Kebangsaan Malaysia, 43600 Bangi, Malaysia

<sup>2</sup> Space Science Centre (ANGKASA), Institute of Climate Change, Universiti Kebangsaan Malaysia, 43600 Bangi, Malaysia

concept of SD is to establish another receiver site at a separation of at least a rain cell at about tens of kilometers (Panagopoulos et al. 2005) to benefit from the inhomogeneity pattern of the rainfall. Both sites are receiving a signal from the same satellite, while the less attenuation sites will be chosen using selection combining or switching technique to be further process at the prime site (Rytir et al. 2017). This concept could be used together with ACM along the way after the selection has been made in case of severe condition of selected signal is still detected (Capsoni et al. 2009).

In tropical region, the number of SD studies is increasing recently. There are three types of data acquisition technique for SD investigation that have been conducted in the literature. One technique is using direct measurement from satellite (Acosta et al. 2012; Cuervo et al. 2016), another technique is using weather radar data (Yeo et al. 2011; Lam et al. 2015) and the third technique is using rainfall data (Islam et al. 2017; Harika et al. 2018). The propagation experiment to investigate the effectiveness of SD scheme in tropical region was first reported by Timothy et al. (2001), where he compared the local gain measurement of two sites separated by 12.3 km at Singapore with ITU-R SD gain model, using lowest claimed baseline orientation angle which was  $4^\circ$ . Another comparison study was conducted by Yeo et al. (2011), also was based in Singapore. He used radar data to compare Hodge and ITU-R gain model and thus concluded that Hodge model was not sensitive to separation distance, while ITU-R model overestimated the gain in tropical region. The experiment was conducted at 18.9 GHz signal with elevation angle of  $44.5^\circ$ . Semire et al. (2014, 2015) investigated on SD link parameters at five different sites based in Malaysia, Indonesia, Philippines, Thailand and Fiji. The author compared Hodge, Panagopoulos and Nagaraja (2012) model at Ku-band signal, thus concluded that Hodge model is practical for separation distance below 10 km and Panagopoulos and Nagaraja models are best suited for temperate region. In year 2017, Islam et al. (2017) experimented SD scheme using rain attenuation deduced from rain intensity measurement. The SD gain models were analyzed using local data, which was at International Islamic University Malaysia and Faculty of Engineering, UKM, Malaysia, separated by 37 km. The comparison was made using elevation angle of  $77.4^\circ$ , frequency 12 GHz and baseline angle of  $0^\circ$  and  $90^\circ$ , between ITU-R, Hodge, Panagopoulos and Semire models (Islam et al. 2017). With this local measurement, ITU-R and Hodge model showed good agreement when using baseline angle of  $90^\circ$ , while Panagopoulos gives better prediction of SD gain when baseline angle of  $0^\circ$  was used. Semire model was observed to underestimate the gain for both  $0^\circ$  and  $90^\circ$  baseline angles, probably because this model was derived from a small-scale data test values of baseline angle.

This study is motivated by the lack of Ka-band SD scheme research in Malaysia and in general in tropical

regions. To the best knowledge of the authors, in year 2015, there was a study on SD scheme using 20.245 GHz frequency at low degree of elevation angle;  $25^\circ$ , focusing on statistical analysis of rain fade dynamics (Jong et al. 2015). Therefore, there are less investigations on SD gain model at Ka band using direct measurement data from satellite with high elevation angle. From the comparison made by researchers in the literature, no prediction model shows consistency with the local data in tropical region. Therefore, the model proposed by scholars of the field particularly focusing on tropical region climate need to be evaluated using various local configurations. In addition, since the distance of rain-cell varies according to local terrain, the optimal exclusion zone for SD scheme should be considered to be as far as possible from the main site to avoid the same coming formation of coming cloud and rain caused by wind blow at nearby locality. SD gain prediction model is significant to assist the telecommunication engineer to estimate the capability of SD, before deploying the facilities. In other words, an effective investment can be performed by satellite operators if researchers can supply them with accurate and reliable tools to measure the level of SD technique, so that they can reduce costs and time to make decisions about the implementation of the scheme (Fenech et al. 2014).

This paper focuses on the SD gain obtained from live measurement of Ka-band signal at two separated location in west Malaysia with high elevation angle of  $68.8^\circ$  and large diameter of gateway antenna. The measured gain was compared with the empirical model developed by Hodge, Panagopoulos, Semire, X. Yeo, and ITU-R model as well. A brief description of each model is narrated at the next section, highlighting the differences of each other. Then, “Methods” is the methods of the measurement, showing that the factors that contribute to the attenuation, and “Results and discussions” discusses the results. Finally, a conclusion is drawn based on the results found.

### Site diversity gain prediction models

Most common metric to measure the SD effectiveness is by calculating the gain and improvement factor. The SD gain is calculated as the difference between the attenuation of single site and attenuation of joint sites at the same percentage of probability of time exceedance. Due to scarce availability of measured data, two types of prediction methods are developed, which are physical and empirical model. This article focuses on empirical models, as it is more easily to be applied and provides faster results than physical models that require a comprehensive understanding of rain process, plus the difficulty of obtaining the data required to be used in the model (Yeo et al. 2015). The first empirical model was proposed in 1981 (Hodge 1981), an initiative from Hodge, that had started the research in early 1970's. It was an improvement from the earlier version

(Hodge 1976), highlighting the fourth factors that contributed to the SD gain, which was the link frequency, with more experimental data than the former. This newly improved model, namely Hodge's model was based on the 34 diversity experiments which was conducted in Canada, England, Japan and the United States, with frequencies ranged from 11.6 to 30 GHz, separation distances from 1.7 to 46.9 km, elevation angle from  $10.7^\circ$  to  $55^\circ$  and baseline orientation angle from  $0^\circ$  to  $164^\circ$  which was then scaled from  $0^\circ$  to  $90^\circ$  (Bosisio et al. 1993; Hodge 1981). This model is a multiplicative of all gains contributed to the total overall site diversity gain,  $G_{SD}$ , that are frequency,  $G_f$ , baseline orientation angle,  $G_\varphi$ , elevation angle,  $G_\theta$  and site separation distance,  $G_d$  as in (1).

$$G_{SD} = G_f G_\varphi G_\theta G_d, \quad (1)$$

$$\begin{aligned} \text{With: } G_d &= a(1 - e^{-bd}), \\ \text{where: } a &= 0.64A - 1.6(1 - e^{-0.11A}), \\ b &= 0.585(1 - e^{-0.098A}), \end{aligned}$$

$$G_f = 1.64e^{-0.025f},$$

$$G_\theta = 0.00492\theta + 0.834,$$

$$G_\varphi = 0.00177\varphi + 0.887.$$

The model was then adopted in ITU-R with a bit modification on the coefficient to suit the ITU-R databank (Ippolito 2017). ITU-R model can be obtained from the guideline document available (ITU-R P.618–13 2017). Hodge model had been compared with another empirical techniques as well namely Goldhirsh, Allnutt and Rogers, and CCIR and physical model, namely Mass, Matricciani and Capsoni et al. model using temperate region databanks (Bosisio et al. 1993). The result was in favor to Hodge model. Panagopoulos et al. (2005) identified that Hodge model was not sensitive to the separation distance and portray gain saturation at a distance more than 15 km. The gain was supposedly to increase proportionally to the distance due to higher decorrelation in rainfall events, which leads to unbalanced attenuation threshold at both sites. Therefore, the author proposed new coefficients to be used in the multiplication of the gain contributors, which is the addition of single site attenuation gain, such that in (2), so that it compensated the justified weaknesses of Hodge model, using temperate region database (Panagopoulos et al. 2005).

$$G_{SD} = G_{A_s} G_d G_f G_\theta G_\varphi \quad (2)$$

$$\text{With: } G_{A_s} = 8.19A_s^{0.0004} + 0.1809A_s - 8.2612,$$

$$G_d = \ln(3.6101d),$$

$$G_f = e^{-0.0006f},$$

$$G_\theta = 1.2347(1 - \theta^{-0.356}),$$

$$G_\varphi = 1 - 0.0006\varphi.$$

Semire et al. (2015) experimented SD using databank from tropical region of five distinct countries, namely Malaysia, Philippines, Indonesia, Thailand and Fiji. Hodge model has been analyzed using the tropical region data and it was found that the model showed less accuracy at lower elevation angle and high frequency. Therefore, the author proposed new expressions and coefficients involving low elevation angles from  $10^\circ$  to  $50^\circ$  and high frequency up to 70 GHz, while the model's structure is unchanged and in line with the Hodge model, such that in (3). This new prediction model was compared with the original Hodge, Panagopoulos and Nagaraja (2012) models and it was apparently deduced that Semire et al. model predicted well than others when tested on data of tropics.

$$G_{SD} = G_d G_f G_\theta G_\varphi, \quad (3)$$

$$\begin{aligned} \text{With: } G_d &= a(1 - e^{-bd}), \\ \text{where: } a &= 0.7755A + 0.3374(1 + e^{-9.16A}), \\ b &= 0.1584(1 + e^{-0.03164A}), \end{aligned}$$

$$G_f = 1.006e^{-0.0015f} - 0.395e^{-0.473f},$$

$$G_\theta = 0.899(1 + \theta^{-0.683}),$$

$$G_\varphi = -0.0000015\varphi + 0.9877.$$

Yeo et al. (2015) derived SD gain prediction model from experimental inference in Yeo et al. (2011) which concluded that the SD gain depends only on three factors; single site attenuation,  $A_s$  site separation distance,  $d$  and elevation angle,  $\theta$ . The model was in different structure than Hodge's model such that in (4) (Yeo et al. 2015). The authors' new model was compared with ITU-R Hodge-based model (empirical model) and ITU-R Paraboni–Barbaliscia (P–B) model (physical model). The result was in favor of Yeo's model which was based on 2025 slant path attenuation of 10–30 GHz frequencies at 45 sites in Singapore with elevation angle ranged from  $10^\circ$  to  $90^\circ$  at intervals of  $20^\circ$ . The site separation distance was varied from 5 to 37 km.

$$G_{SD} = (-0.78 + 0.88A_s)(1 - e^{-0.18d})(1 + e^{-0.14\theta}). \quad (4)$$

## Methods

A live measurement to monitor the SD scheme has been conducted at gateway stations, Cyberjaya ( $101.6584^\circ$  E,  $2.9356^\circ$  N) and Rawang ( $101^\circ 33' 16.6''$  E,  $3^\circ 18' 13.1''$  N) separated by a direct distance of 42.52 km. This monitoring activities are currently running by MEASAT Satellite System Sdn. Bhd.,

at Ka-band frequency (20.2 GHz) using large diameter earth station antenna of 7.3-m with antenna gain of 62.31 dB, elevation angle of 68.8° and vertical polarization. The antenna received signals from MEASAT-5, a high throughput satellite (HTS), with 87 Ku-band and 10 Ka-band transponders, deployed by space Systems/Loral (SS/L) contractors, located at 119.5° E with 63dBW EIRP. The beacon receiver samples in time series with 2 s interval at Cyberjaya and 1 s interval at Rawang, stored in daily basis. The slant path attenuation at both locations was obtained by averaging each data to 1-min, thus deducting it with the clear sky value of the day. The spikes and outliers that were found in the time series data were eliminated, and 2 days (four files with the same date at Cyberjaya and Rawang) that consisted of data anomalies were removed as well. The anomalies might be caused by system malfunction or system under maintenance as they appeared to begin at the same point of time, such that in Fig. 1. Cyberjaya showed longer period of data unavailability in Fig. 1. The percentage of data availability in average was 95.13% at Cyberjaya and 95.4% at Rawang for four years of measurements from January 1st, 2014 to December 31st, 2017. Further analyses and single site results could be viewed in Samat and Mandeep (2020).

The averaged 1-min CCDF attenuation graphs were to allow the determination of joint attenuation. An internal program was developed to screen the data to get the lowest attenuations amongst two sites. The program compares every 1-min average data for each day in Cyberjaya with the same data and day in Rawang. Two rain gauge was located near the stations, to record the rainfall in 1 min. The rain gauge was installed by Department of Drainage and Irrigation (DID) Selangor, Malaysia to monitor the rain pattern at Cyberjaya and Rawang area, from January 1st, 2014 to December 31st, 2017. Rain value in millimeter (mm) was converted to millimeter per hour (mm/h) to represent the rain rate by multiplying each with 60. The 1-min rainfall was arranged in ascending order to calculate the frequencies of the same rain rate. Then, the probability of occurrences was derived using the accumulated frequencies divided by total

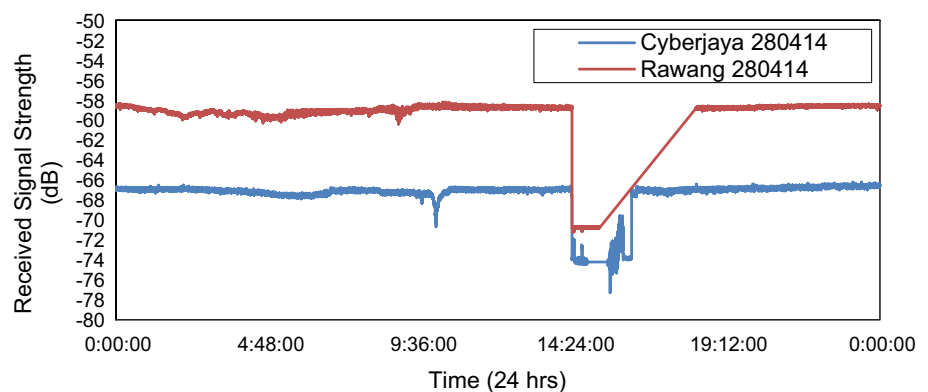
amount of minutes for 1 month, thus the rain rate at percentage of probability of interest could be determined.

Baseline orientation angle was determined by measuring the angle between azimuth line and baseline distance of Cyberjaya and Rawang, as shown in Fig. 2. The baseline orientation angle was carefully ascertained based on the guideline given in the literature (Hodge 1981; Ippolito 2017; Panagopoulos et al. 2005). Therefore, the significant parameters involved in this measurement are as listed in Table 1. From Table 1, the parameter associated with the SD prediction model was taken as input. Each SD gain prediction model was reconstructed using Microsoft EXCELL including ITU-R P.618–13 diversity gain model. The average of 4 years measurement slant path attenuation at both sites was determined and was taken as input to the model as well. All results were discussed in the next section.

## Results and discussions

From the experiment, the attenuation values were plotted and the percentage of time exceedance at 1%, 0.1% and 0.01% were observed. Figures 3 and 4 shows the average monthly attenuation graph at Rawang and Cyberjaya for 4 years started from 1st January of 2014 to 31st December of 2017. From Fig. 3, the highest average of attenuation was observed in April and November in Rawang. November is within northeast monsoon which started from November to March, and inter-monsoon is between end of March to end of April, then May to September is the southwest monsoon in Malaysia's weather (Omotosho et al. 2017). The month of February experienced the lowest attenuation among all the months in the 4 years. From Fig. 3, the average attenuation for percentage of signal unavailability at 1% of time was around 3 dB, 0.1% was in between 4 and 8 dB and 0.01% was around 10–30 dB, then the twelve's graphs tend to saturate above the 30 dB value. This might be due to the limitation of dynamic range value of the beacon receiver. At the same time, this beacon also received high noise from the

**Fig. 1** Data anomalies on 28th of April 2014 during the same period at Cyberjaya and Rawang.



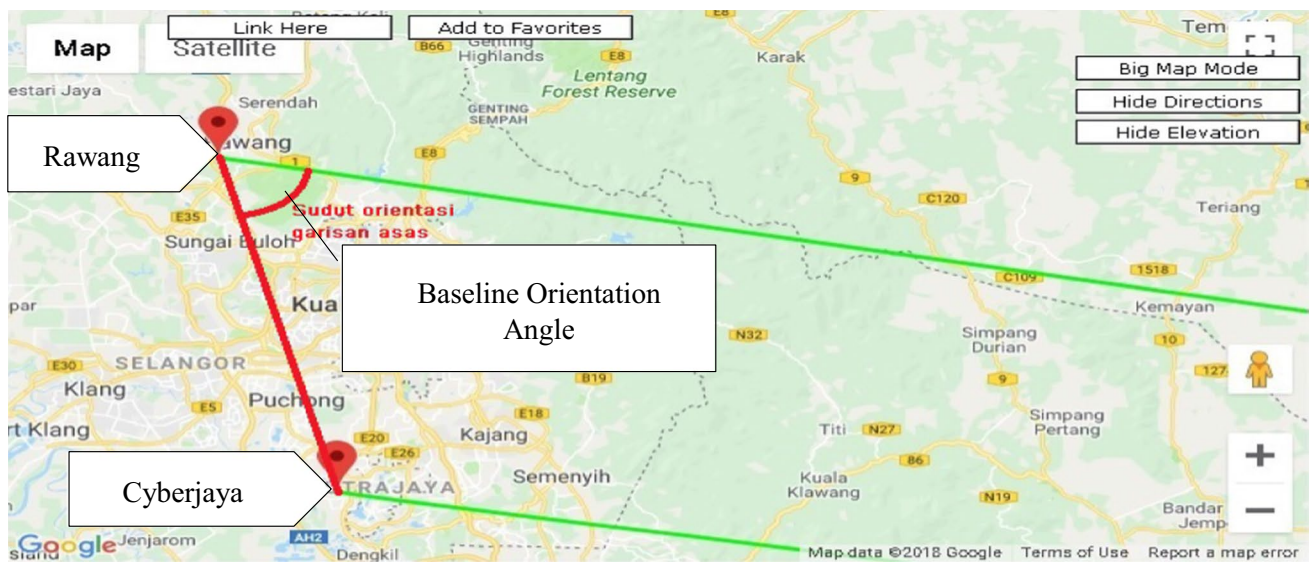


Fig. 2 Baseline orientation angle, an angle between azimuth line (green line) and baseline direct distance of Cyberjaya and Rawang

Table 1 Parameter for SD gain measurement

Parameter	Cyberjaya	Rawang
Frequency (GHz)	20.2	20.2
Azimuth angle	99.3°	100.4°
Altitude (km)	0.01962	0.038
Elevation angle (°)	68.8°	68.8°
Antenna diameter (m)	8.1	8.1
Polarization	Vertical	Vertical
Baseline orientation angle	65°	65°

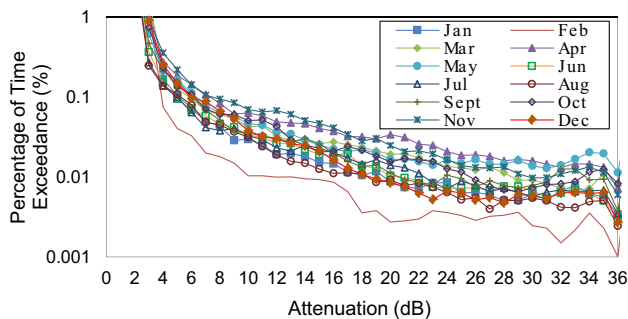


Fig. 3 Average monthly attenuation from year 2014 to 2017 in Rawang

system as well as airspace causing it to mistakenly identify the noise as a signal from the satellite.

The graph of attenuation in Cyberjaya from 2014 to 2017 in Fig. 4 depicted that the attenuation range of over the time of 1% was 1.5–2 dB, 0.1% was 2.6–5 dB and at 0.01% was from 8.2 to 24 dB, then the graph saturated and flat at

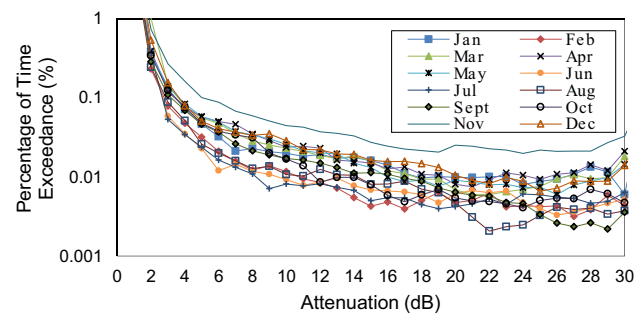


Fig. 4 Average monthly attenuation from year 2014 to 2017 in Cyberjaya

around 26–27 dB. From this whole graph, it is observed that November was experiencing greater attenuation over the other months, as seen from the state of the graph being separated between the other graphs. The least attenuation was experienced in July when considering at 0.01% of signal outage time. Afterall, the attenuation value experienced in Cyberjaya at all months was less than Rawang. This condition is seen in relation to Fig. 6 which shows the low average rain intensity at 0.01 percentage of time exceedance of that site for 4 years compared to Rawang as in Fig. 5.

From Fig. 5, the rain rate was high in November at Rawang and the least was in February, the same goes at Cyberjaya in Fig. 6. The value of this monthly average accumulation is consistent due to northeast monsoon which is occurred in November to March of every year. The average of rain rate at 0.01 percent of time at Rawang ranges from 45 to 116 mm/h, and Cyberjaya ranges from 39 to 112 mm/h. The average rain intensity for each year



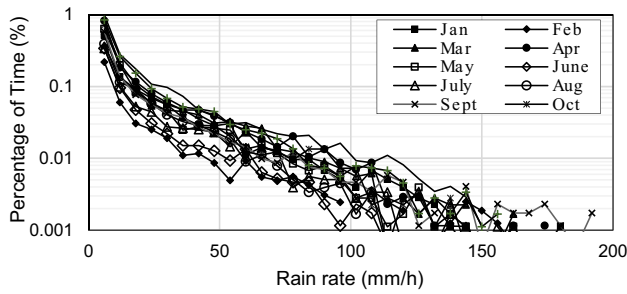


Fig. 5 Average (2014–2017) rain rate of each months at Rawang

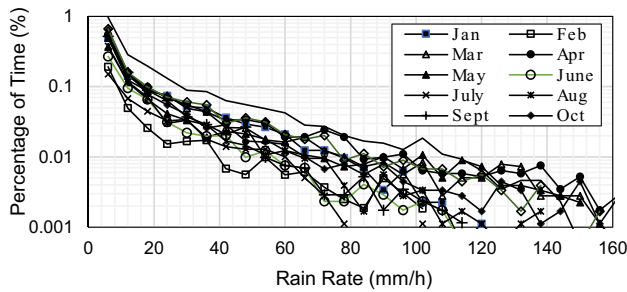


Fig. 6 Average (2014–2017) rain rate of each months at Cyberjaya

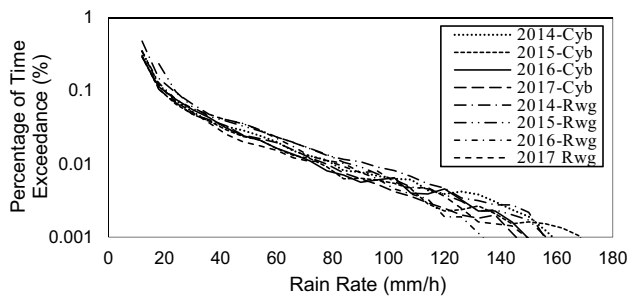


Fig. 7 Average rain intensity from year of 2014 to 2017 at Cyberjaya (Cyb) and Rawang (Rwg)

is displayed in Fig. 7. The rain intensity of 0.01 percentage of time exceedance at Cyberjaya was 80 mm/h, 76 mm/h, 74 mm/h and 78 mm/h in year of 2014, 2015, 2016 and 2017, respectively. While at Rawang, rainfall was more intense than Cyberjaya such that 92 mm/h, 82 mm/h, 84 mm/h and 81 mm/h in year 2014, 2015, 2016 and 2017, respectively, at the same percentage of time. Even though statistics conducted by Shayea et al. (2018) presented the average of 120 mm/h, the collected data were from June 2011 to May 2012. The author also admitted that the ITU-R latest prediction rain rate for Malaysia is still under investigation. Nevertheless, the latest ITU-R P.837–7 stated that the rain intensity in Malaysia is about

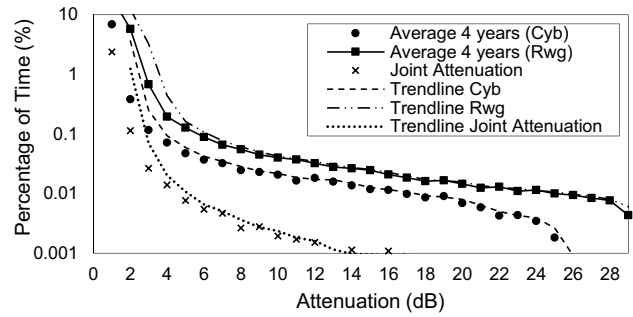


Fig. 8 Average 4 years measurement of attenuation CCDF at Cyberjaya and Rawang and their joint attenuation.

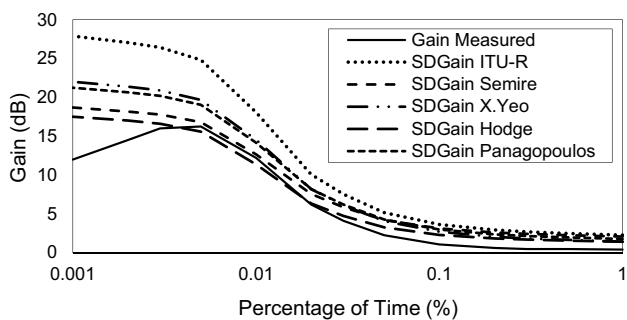
100 mm/h. Therefore, in this case, the measurement of rain data did not differ much than predicted by ITU-R.

The average of attenuation and joint attenuation of year 2014–2017 at each site was deduced and compared. Figure 8 depicted that the trendline (average of two period) attenuation at Rawang for 4 years measurement was 25.8 dB and at Cyberjaya was 17.4 dB. The differences between joint attenuation and single site attenuation at the same percentage of time were noted as gain. The site diversity improvement factor (IF) is calculated by taking the ratio of percentage of time of the same attenuation point at both sites, such that in (5).  $P_s$  is the percentage of time exceedance of single site (main sites), whereas  $P_d$  is the percentage of time exceedance of the joint graph at the same attenuation value. In this case, taking 12 dB of attenuation in Fig. 8 correlates to 0.02% of time exceedance at the Cyberjaya site and about 0.0016% of time exceedance at the diversity graph. Therefore, at this point of attenuation, the IF was 12.5, which mean a great improvement. This leads to meaning that for the same attenuation threshold, the signal unavailability could be improved to 0.0016% of time of an average year using site diversity scheme instead of only 0.02%.

$$IF = \frac{P_s(A)}{P_d(A)} \tag{5}$$

From Fig. 8, it was measured that the gain obtained at 0.01 percentage of time was 12.2 dB, measured between joint attenuation and Cyberjaya (main) site. The induced gain was compared with site diversity gain models as in Fig. 9.

All site diversity gain models deviated far from the measured gain at 0.1 percentage of time, as shown in Fig. 9. When percentage of time exceedance was approaching 0.02%, the Hodge model predicted gain was similar with the measured gain, then it goes a bit underestimated the gain at 0.01% of time. At this point of percentage of time, Semire model apparently predicted the closest value as the measured gain. Following through the decreasing of time percentage, it can be



**Fig. 9** Comparison of each model; ITU-R, Semire, X.Yeo, Hodge and Panagopoulos with the measured gain

**Table 2** RMSE value of each model

Models	RMSE
ITU-R	0.578725742
Semire	0.507469426
X.Yeo	0.425821738
Hodge	0.383566864
Panagopoulos	0.476878661

noted that the measured gain was saturated at 0.005% of the time as what can be observed from the source of the original attenuation. Up to this point of percentage of time that the gain saturated, Semire model seems to have equal consistency of predicted gain value to the measured data. To further analyze the performance of each model, a statistical evaluation was performed according to ITU-R P.311–13 (2009). The root mean square error (r.m.s.e) was identified using formula (7). The test variable  $T_i$  was obtained from the logarithm of the ratio of predicted gain,  $G_p$  and measured gain,  $G_m$ , formulation of (6). For the measured gain less than 10 dB, a scaling factor is applied and without scaling factor for gain value of greater than 10 dB. The r.m.s.e values, denoted as *rmse*, was derived from the calculated mean,  $\mu_T$  and its deviation,  $\sigma_T$  for each percentage of time. While  $N$  is the number of test variables and  $i$  is the count variable up to  $N$ . Table 2 shows the r.m.s.e values of each models.

$$S_i = \frac{G_{pi}}{G_{mi}}, \tag{6}$$

$$T_i = \left\{ \begin{array}{l} (\frac{G_{mi}}{10})^{0.2} \times \ln(S_i) \text{ for } G_{mi} < 10\text{dB} \\ \ln(S_i) \text{ for } G_{mi} \geq 10\text{dB} \end{array} \right\}, \tag{7}$$

$$\mu_T = \frac{1}{N} \sum_{i=1}^N T_i, \tag{8}$$

$$\sigma_T = \sqrt{\frac{1}{N} \sum_{i=1}^N (T_i - \mu_T)^2}, \tag{9}$$

$$\text{Withrmse} = \sqrt{\mu_T^2 + \sigma_T^2}. \tag{10}$$

From the results of r.m.s.e, it is found that the Hodge model predicted the smallest error as a whole, calculated from a time exceedance of 1–0.005%, that is, before the time the measured graph gain changes to fall down. However, Hodge model tends to underestimate the gain from 0.005% and downward if the measured graph were to be extrapolated up to 0.001% of time. It was obvious that the second least error was X. Yeo model, followed by Panagopoulos, Semire and ITU-R model. From the analysis of the models’ shape, X. Yeo and Panagopoulos model were portrayed similar shape with the measurement graph, which shows abrupt negative slope as the percentage of time decreases, at point 0.02% of time and downward. However, the gain predicted by both models is higher than the measured one, as if the measured gain was positively shifted at y-axis at a certain scale. X. Yeo model seems the most suitable prediction model for this sample data, since it predicted less error than Panagopoulos model, and another merit is that it is having the same shape pattern of graph with the measured data. While Semire model apparently has the least error from 0.02% of time up to 0.005% of time, the model predicted bigger error from 1 to 0.02% of time percentage. This contributes to larger r.m.s.e to Semire model’s prediction.

### Conclusion

The measurement of this study was obtained from two sites of Malaysia, at Cyberjaya and Rawang, of Ka-band HTS signal with large diameter of antenna and high elevation angle. The gain at 0.01% of time was measured as 12.2 dB, and when comparing with other models, at the same percentage of time, Semire Model shows the most suitable for the measured data. However, in overall, from the percentage of time exceedance of 1% up to 0.005%, where the point of attenuation saturated, none of models shows a comprehensive match to the measured gain. Though Hodge model shows the least r.m.s.e., yet from the visual observation of the graph, the lowest r.m.s.e value does not mean that it successfully predicted equal value as the measured gain for the whole percentage of time exceedance, only it portrayed less error than other models. The closest pattern of models’ graph with the measured one was X. Yeo and Panagopoulos models, but both models overpredicted the gain. More research should be explored to further digging the characteristics and factors that contribute to the development of site

diversity gain model especially that is related to the pattern differences of two site's attenuation thresholds, which are mainly influenced by the intensity of the local rain. Therefore, from the observation of the analysis, the gain contribution to develop a prediction model such as frequency, base-line angle, elevation angle and separation distance should be analysed further.

**Acknowledgements** We thank MEASAT Satellite Systems Sdn Bhd, a leading satellite operator in Malaysia for providing us with the data from its earth station gateways in Cyberjaya and Rawang for the years of 2014 to 2017. The authors also would like to appreciate and thank Mr. Ahmad Zaidee bin Abu, an independent researcher, for his contributions in developing scripts to process data according to the given algorithm.

## Compliance with ethical standards

**Conflict of interests** On behalf of all authors, the corresponding author states that there is no conflict of interest.

## References

- Acosta R, Morse J, Zemba M, Nessel J (2012) Two years of site diversity measurements in Guam, USA. In: 18th Ka and Broadband Communications, Navigation and Earth Observation Conference, 24–27 September 2012, Ottawa, Canada, available from: <https://ntrs.nasa.gov/archive/nasa/casi.ntrs.nasa.gov/20120016399.pdf>. Accessed 22 Jun 2019.
- Bosisio AV, Capsoni C, Matricciani E (1993) Comparison among prediction methods of site diversity system performances. In: Eighth International Conference on Antennas and Propagation, 30 March–2 April 1993, Edinburgh, UK, 60–63, ISBN: 0–85296–572–9.
- Callaghan SA, Boyes B, Couchman A, Waight J, Walden CJ, Ventouras S (2008) An investigation of site diversity and comparison with ITU-R Recommendations. *Radio Sci* 43(RS4010):1–8. <https://doi.org/10.1029/2007RS003793>
- Capsoni C, D'Amico M, Nebuloni R (2009) Time and site diversity gain: a close relationship. In: International Workshop on Satellite and Space Communications, 9–11 September 2009, Tuscany, Italy, 166–170, DOI: 10.1109/IWSSC.2009.5286387.
- Cuervo F, Schönhuber M, Capsoni C, Hong Yin L, Jong SL, Bin Din J, Martellucci A (2016) Ka-Band propagation campaign in Malaysia—first months of operation and site diversity analysis. In: 10th European Conference on Antennas and Propagation (EuCAP), 10–15 April 2016, Davos, Switzerland, 1–5, DOI: 10.1109/EuCAP.2016.7481248.
- Fenech H, Tomatis A, Amos S, Serrano Merino J, Soumpholphakdy V (2014) An operator's perspective on propagation. In: 8th European Conference on Antennas and Propagation (EuCAP), 6–11 April 2014, The Hague, Netherlands, 3349–3352, DOI: 10.1109/EuCAP.2014.6902545.
- Harika S, Nagarjuna S, Naveen TV, Sanjay Harshanth G, Sri Kavya KC, Kotamraju SK (2018) Analysis of rain fade mitigation using site diversity technique in southern tropical region of India. *Int J Eng Technol* 7(1):622–626. <https://doi.org/10.14419/ijet.v7i1.1.10790>
- Hodge DB (1976) An empirical relationship for path diversity gain. *IEEE Trans Antennas Propag* 24(2):250–251. <https://doi.org/10.1109/TAP.1976.1141304>
- Hodge DB (1981) An improved model for diversity gain on earth-space propagation paths. *Radio Sci* 17(6):1393–1399. <https://doi.org/10.1029/RS017i006p01393>
- Ippolito LJ Jr (2017) *Satellite communications system engineering. Atmospheric effects, satellite link design and system performance*, 2nd edn. John Wiley and Sons Ltd, West Sussex
- Islam MR, Habaebi MH, Haidar IMD, Lwas AK, Zyoud A, Mandeep S (2015) Rain fade mitigation on earth-to-satellite microwave links using Site Diversity. In: IEEE 12th Malaysia International Conference on Communications (MICC), 23–25 November 2015, Kuching, Malaysia, 186–191, DOI: 10.1109/MICC.2015.7725431.
- Islam MR, Lwas AK, Habaebi MH (2017) Site diversity gain for earth-to-satellite links using rain intensity measurement. *Indones J Electric Eng Inf (IJEEI)* 5(4):330–338. <https://doi.org/10.11591/ijeei.v5i4.364>
- Jong SL, Lam HY, Din J, Amico MD (2015) Investigation of Ka-band satellite communication propagation in equatorial regions. *ARNP J Eng Appl Sci* 10(20): 9795–9799, available from: <https://core.ac.uk/download/pdf/42955742.pdf>. Accessed 23 Jun 2019.
- Kyrgiazos A, Evans B, Thompson P (2014) On the gateway diversity for high throughput broadband satellite systems. *IEEE Trans Wirel Commun* 13(10):5411–5426. <https://doi.org/10.1109/TWC.2014.2339217>
- Lam HY, Luini L, Din J, Capsoni C, Panagopoulos AD (2015) Performance of site-diversity satellite communication systems in equatorial Malaysia investigated through weather radar data. In: 9th IEEE European Conference on Antennas and Propagation (EuCAP), 13–17 April 2015, Lisbon, Portugal, 1–4, available from: <https://ieeexplore.ieee.org/document/7228442>. Accessed 23 Jun 2019.
- Nagaraja C, Otung IE (2012) Statistical prediction of site diversity gain on earth-space paths based on radar measurements on the UK. *IEEE Trans Antenna Propag* 60(1):247–256. <https://doi.org/10.1109/TAP.2011.2167896>
- Omotosho TV, Mandeep JS, Abdullah M (2011) Cloud-cover statistics and cloud attenuation at Ka- and V-Bands for satellite systems design in tropical wet climate. *IEEE Antennas Wirel Propag Lett* 10:1194–1196. <https://doi.org/10.1109/LAWP.2011.2172674>
- Omotosho TV, Akinwumi SA, Ometan OO, Adewusi MO, Mandeep JS, Abdullah M (2017) Earth-Space rain attenuation prediction: its impact at Ku, Ka and V-band over some equatorial stations. *J Inf Math Sci* 9(2):359–374
- Panagopoulos AD, Arapoglou PDM, Cottis PG (2004) Satellite communications at KU, KA, and V bands: Propagation impairments and mitigation techniques. *IEEE Commun Surv Tut* 6(3):2–14. <https://doi.org/10.1109/COMST.2004.5342290>
- Panagopoulos AD, Arapoglou PDM, Kanellopoulos JD, Cottis PG (2005) Long-term rain attenuation probability and site diversity gain prediction formulas. *IEEE Trans Antennas Propag* 53(7):2307–2313. <https://doi.org/10.1109/TAP.2005.850762>
- ITU-R Recommendation P.311–13 (10/2009). Acquisition, presentation and analysis of data in studies of tropospheric propagation. Geneva. Electronic Publication.
- ITU-R Recommendation P.618–13 (2017) Propagation data and prediction methods required for the design of earth-space telecommunication systems. Geneva. Electronic Publication.
- Rytir M, Cheffena M, Grotthing PA, Braten LE, Tjelta T (2017) Three-site diversity at Ka-band satellite links in Norway: gain, fade duration and the impact of switching schemes. *IEEE Trans Antennas Propag* 65(11):5992–6001. <https://doi.org/10.1109/TAP.2017.2751667>
- Samat F, Mandeep JS (2019) Rain attenuation at tropical region site diversity gain models sensitivity. *Indone J Electric Eng Inf (IJEEI)* 7(3):472–483. <https://doi.org/10.11591/ijeei.v7i3.956>
- Samat F, Mandeep JS (2020) Impact of rain attenuation to Ka-Band signal propagation in tropical region: a study of 5-Year MEASAT-5's

- beacon measurement data. *Wirel Personal Commun.* <https://doi.org/10.1007/s11277-020-07172-x>
- Semire FA, Rosmiwati M, Widad I, Norizah M, Mandeep JS (2014) Evaluation of site diversity rain attenuation mitigation technique in South-East Asia. *J Acta Astronautica Sci Direct* 96(1):303–312. <https://doi.org/10.1016/j.actaastro.2013.11.034>
- Semire FA, Mohd-Mokhtar R, Ismail W, Mohamad N, Mandeep JS (2015) Modeling of rain attenuation and site diversity predictions for tropical regions. *J Ann Geophys* 33(3):321–331. <https://doi.org/10.5194/angeo-33-321-2015>
- Shayea I, Rahman TA, Hadriazmi M, Islam MR (2018) Real measurement study for rain rate and rain attenuation conducted over 26 GHz microwave 5G link system in Malaysia. *IEEE Access* 6:19044–19064. <https://doi.org/10.1109/ACCESS.2018.2810855>
- Timothy KI, Ong JT, Choo EBL (2001) Performance of the site diversity technique in Singapore: preliminary results. *IEEE Commun Lett* 5(2):49–51. <https://doi.org/10.1109/4234.905932>
- Yang R, Li L, Zhao Z, Lu T (2013) Cloud simulation and attenuation at Ka band on slant path. In: 2013 Cross Strait Quad-Regional Radio Science and Wireless Technology Conference, 21–25 July 2013, Chengdu, China. DOI: 10.1109/CSQRWC.2013.6657414.
- Yeo JX, Lee YH, Ong JT (2011) Performance of site diversity investigated through RADAR derived results. *IEEE Trans Antennas Propag* 59(10):3890–3898. <https://doi.org/10.1109/TAP.2011.2163770>
- Yeo JX, Lee YH, Ong JT (2015) Site diversity gain at the equator: radar-derived results and modelling in Singapore. *Int J Satell Commun Netw* 33(2):107–118. <https://doi.org/10.1002/sat.1074>
- Yuan F, Lee YH, Meng YS (2016) Investigation of cloud attenuation on Ka-band satellite beacon signal in tropical region. In: 2016 IEEE International Symposium on Antennas and Propagation (APSURSI), 26 June–1 July 2016, Fajardo, Puerto Rico. DOI: 10.1109/APS.2016.7696335.
- Yuan F, Lee YH, Meng YS, Yeo JX, Ong JT (2017) Statistical study of cloud attenuation on Ka-band satellite Signal in Tropical Region. *IEEE Antennas Wirel Propag Lett* 16:2018–2021. <https://doi.org/10.1109/LAWP.2017.2693423>
- Yussuff AIO, Hamzat N, Khamis NHH (2017) Site diversity technique application on rain attenuation for Lagos. *Indones J Electric Eng Inf (IJEEI)*. 5(1):77–84. <https://doi.org/10.11591/ijeei.v5i1.262>



# Performance analysis of IRNSS using compact microstrip patch antenna for S band application

Mitchell Prajapati<sup>1</sup> · Abhishek Rawat<sup>1</sup>

Received: 11 July 2019 / Accepted: 24 June 2020 / Published online: 7 July 2020  
© Institute of Geophysics, Polish Academy of Sciences & Polish Academy of Sciences 2020

## Abstract

The Indian Regional Navigation Satellite System (IRNSS) is a newly functional regional satellite navigation system around Indian subcontinent. It is also known as its functioning name of NaVIC, Navigation with Indian Constellation, working on L5 and S band. In this paper a compact and low cost circularly polarized microstrip patch antenna is proposed for IRNSS S band application, to provide single frequency navigation solution. Fabrication of proposed antenna is performed using substrate of RT Duroid 5880 with the dimensions of  $0.66\lambda \times 0.5\lambda$ . The performance of IRNSS is investigated by prototype antenna and Accord made triband antenna with IRNSS/GPS/SBAS receiver which is enable to receive L1, L5 and S band data. The comparison of proposed S band antenna, triband accord system antenna and GPS antenna is presented in terms of carrier to noise ratio and positioning error. Results show that proposed antenna is suitable for S band application of IRNSS receiver.

**Keywords** Indian Regional Navigation Satellite System · Microstrip patch antenna · Duroid · PVT-position, velocity, timing · IRNSS receiver · GPS

## Introduction

Indian Regional Navigation Satellite System (IRNSS) is established and controlled by the Indian Space Research Organization (ISRO) under Government of India. This independent navigation system was required for India since long back, because other global positioning and navigation systems are not reliable in inimical conditions. The main objective of the IRNSS is to provide positioning and navigation services to users in the Indian region. IRNSS is designed to provide navigation solutions to land, air and marine transport users (Ganeshan et al. 2005) in place of GPS. IRNSS provides approximately 10-m accuracy for positioning in the Indian region and 20 m accuracy for positioning in 1500 km around the Indian region. There are seven satellites to make the IRNSS fully operational since May 2016. IRNSS transmits signals in dual band namely L5 with centre frequency

of 1176.45 MHz and S band with centre frequency of 2492.028 MHz. So dual bands as well as any of the band can be used for the navigation purpose (ISRO 2019). Since S band is less prone to ionosphere effect, we focus on the S band antenna design. We also compare it with Accord made triband antenna and GPS antenna in terms of signal strength, carrier to noise ratio with positioning accuracy of IRNSS/GPS/SBAS receiver.

Nowadays, very portable handheld wireless communication devices are in high demand for internet and mobile communication. Antenna size and its gain is the major consideration for these devices. Compact microstrip patch antenna is a better choice with portable communication devices. Navigation is one of the important application of satellite communication. For radio navigation, highly accurate receiver unit is required with its high gain antenna performance due to the longer distance between satellite and ground receiver (Wu et al. 2010; Bilotti and Vegni 2010). Circular polarization is preferred in satellite communication to overcome the limitation of transmit and receive antenna orientation. (Sahal and Tiwari 2016) The axial ratio for circular polarization must be within the 3 dB to achieve circular polarization. Microstrip patch antenna is a thin flat structured antenna in which some simple techniques can be applied to get circular polarization, so it is preferable for satellite receivers.

✉ Mitchell Prajapati  
mishelprajapati@gmail.com

Abhishek Rawat  
arawat@iitram.ac.in

<sup>1</sup> Institute of Infrastructure Technology Research and Management, Near Khokhara Circle, Maninagar East, Ahmedabad, Gujarat 380026, India

In this paper a low cost compact design of microstrip patch antenna is proposed for S band IRNSS application with circular polarization. A prototype of proposed antenna is fabricated and tested. Real time navigation signal, received with proposed antenna is closed confirmation with signal received by Accord made triband antenna. The positioning information and its error are found and discussed with proposed antenna, Accord made triband antenna and GPS antenna. Here, second section includes proposed “S band antenna design”, third section contains “Experimental setup”, fourth section comprises “Results and discussion”, and fifth section concludes the findings.

## S band antenna design

The proposed design consists of conducting patch and ground plane, made up of copper with substrate of duroid with dielectric constant of 2.2. It contains loss tangent of 0.0009 which exhibits very low dielectric loss. The antenna is fed with coaxial cable. The dimensions of the patch are 38 mm × 29 mm, tuned to achieve 2.49 GHz in S band. Dimensions of the substrate are 80 mm × 60 mm with thickness of 3.2 mm to get desired return loss and bandwidth. There are several methods to achieve circular polarization in microstrip patch antennas, like cross slits, truncated corners, and dual feed excited by two orthogonal modes technique etc. (Kumar and Ray 2003). The truncated-corners microstrip patch antenna is the best choice for small axial ratio with narrow axial ratio bandwidth (Sharma and Gupta 1983; Sahal and Tiwari 2016). Cross shaped slits gives circular polarization with good axial ratio bandwidth (Nasimuddin and Qing 2012). Feed location and feeding technique decides the impedance bandwidth and axial ratio which is ratio of minor and major axis of polarization circle or ellipse and it must be less than 3 dB, decides the polarization (Sahal and Tiwari 2016). As the circular polarization is preferred for satellite communication, half circular cuts are created at both diagonal corners to set the RHCP axial ratio is around 3 dB.

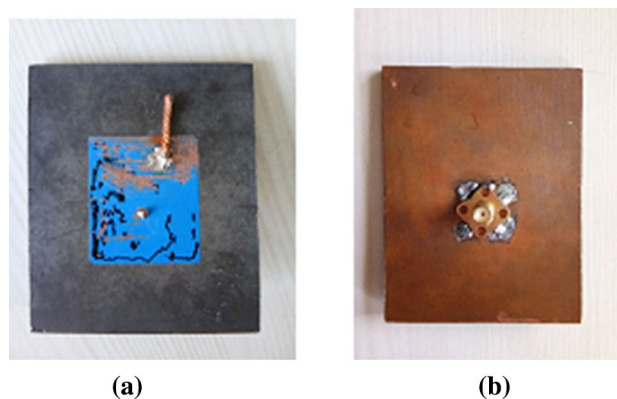
The proposed design of circularly polarized Microstrip patch antenna is focused for S band of IRNSS receiver with the resonant frequency of 2.4900 GHz. This design are prepared and simulated in Ansoft HFSS 17. The prototype of proposed design is also prepared and initially tested on Keysight Field fox Microwave Analyzer N9916A.

Return loss is the loss of power which is returned back to the antenna. As much as the value of return loss is small, the antenna can radiate more in forward direction. So the smaller (negative) value of return loss is preferred. Generally, the return loss must be less than -10 dB is considered for practical application. For proposed antenna, the return loss found to be of -27.85 dB with the bandwidth of 60 MHz. Impedance bandwidth for simulated antenna is 2% and prototype

antenna is 2.41%. The fabricated antenna received the resonant frequency of 2.5422 GHz which was higher than the desired frequency. The simulation results of the design are not matched with prototype results when tested on microwave analyzer and slight frequency deviation is found. So, the tuning stub is used on the edge of patch, to reduce the resonant frequency.

This stub will increase the effective length of patch and shift the axis of field electrically, also the position of feed point will be changed electrically. The effective length will determine the resonant frequency and effective position of feed point will decide the input impedance (Reddy et al. 2015), resulting the desired frequency is achieved with very good value of return loss. When the length of the stub is very small, less than  $\lambda/4$ , then by changing its length and width, the resonant frequency of the microstrip patch antenna is tuned (Ray and Kumar 2000). Generally, the 10% tuning range of frequency is achieved by changing the length of stub from 0 to 1 cm (Roy and Jha 2019). The copper stub (as shown in Fig. 1a) of 19 mm × 2 mm is attached along the length of the patch which helped to set the desired resonant frequency. This antenna radiates on the center frequency of 2.4845 GHz with the bandwidth of 60 MHz which covers S band operating bandwidth of IRNSS. The gain is 7.43 dBi and standing wave ratio of 1.24 dB. All these operating parameters make our antenna fully compatible to IRNSS application.

A planar microstrip patch antenna is having dielectric substrate between conducting ground plane and patch, due to which fringing effect is generated. It makes the effective dielectric constant always less than relative dielectric constant of substrate. So effective width and effective length of patch is considered for further calculations. The effect of stub increases the overall resonant length of microstrip patch antenna from  $L_e$  to  $L_e + \Delta L_1$ , the value of  $\Delta L_1$  can be found by following equation (Kumar and Ray 2003)



**Fig. 1** **a** Top view of fabricated antenna, **b** bottom view of fabricated antenna

$$\Delta l_1 = \frac{(w_e)(l_e)}{W_e} \tag{1}$$

where  $(w_e)(l_e)$  is effective area of stub and  $W_e$  is effective width of patch.

Now the new resonant frequency  $f_0$  will be (Kumar and Ray 2003)

$$f_0 = \frac{c}{2(L_e + \Delta l_1)} \sqrt{E_{\text{eff}}} \tag{2}$$

where  $c$  is velocity of light in meter per second,  $L_e$  is effective length for patch and  $E_{\text{eff}}$  is effective dielectric constant of substrate.

As the new desired frequency is known to us, we can find out the approximated length and width of stub by putting the value of  $f_0$  in Eq. (2) (Kumar and Ray 2003) (Table 1).

Here, Fig. 1a shows photograph of front side of fabricated antenna with corner truncated patch on the duroid substrate and a piece of copper wire is soldered on the surface of the patch for stub matching. Figure 1b shows the photograph of

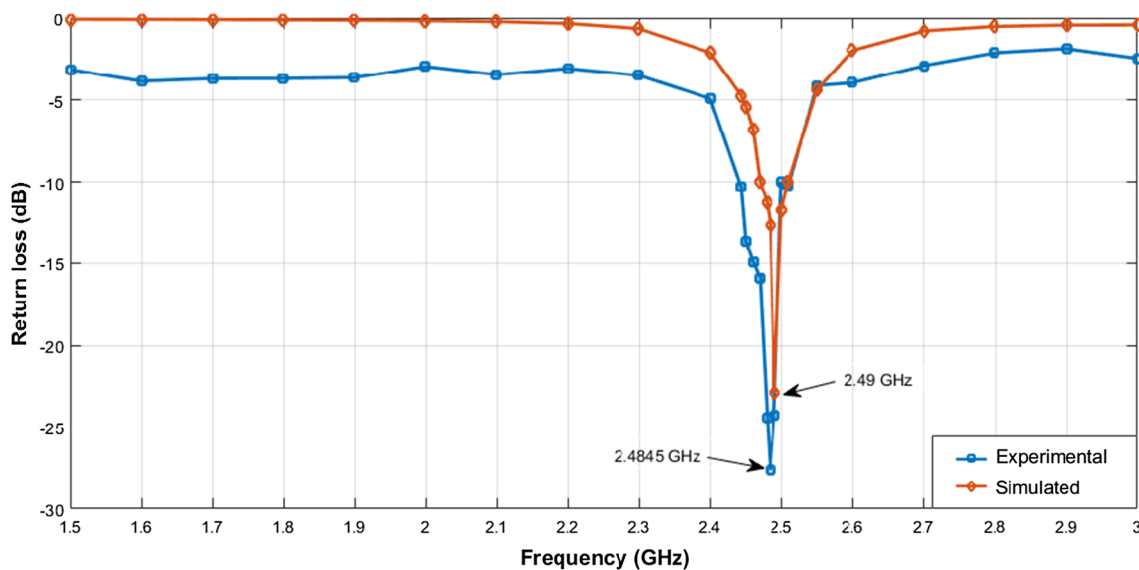
**Table 1** Comparison of simulated and measured results of prototype antenna

Parameters	Simulated	Measured
Resonant frequency	2.4900 GHz	2.4845 GHz
Return loss	− 22.9331 dB	− 27.85 dB
Gain	7.89 dB	7.43dBi
Bandwidth	52 MHz	60 MHz
Impedance bandwidth in %	2%	2.41%

back side of the fabricated antenna, over which the SMA connector is soldered to feed the antenna. Figure 2 shows the simulated and measured values of return loss (reflection coefficient), S11 in dB of proposed antenna. The red curve shows the simulated data and blue curve shows the experimental data of return loss. Table 2 shows the comparison of antenna dimensions, patch type and size, bandwidth and gain of prototype antenna with other antennas present in the literature, and found that the size of proposed antenna is smaller with its gain value, comparing with other antennas.

### Experimental setup

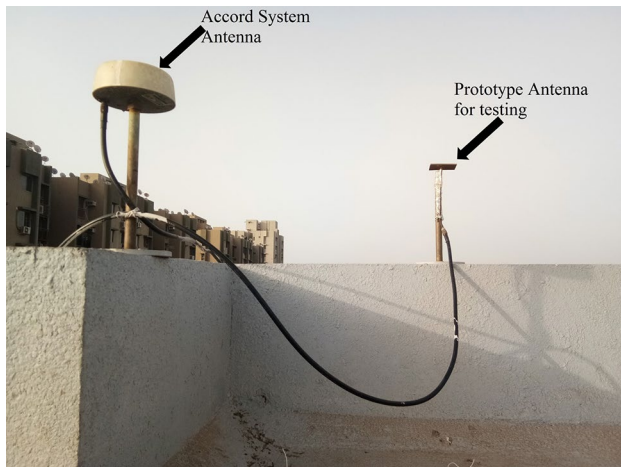
Ansys HFSS 17.2 is used to design and simulate the circularly polarized microstrip patch antenna for the resonant frequency of 2.49 GHz with RHCP axial ratio around 3 dB. The proposed design is fabricated using duroid substrate. Prototype of the antenna is initially tested with Keysight Field fox Microwave Analyzer N9916A and found that the results deviate from desired band. By properly designing and placing the stub, along the length of the patch, the resonant frequency of 2.4845 GHz is achieved with the return loss of − 27.85 dB for S band operation. Afterward this prototype antenna is connected with the IRNSS/GPS/SBAS receiver to receive the S band navigation signal by setting the receiver into S band mode only. There are two sets of IRNSS/GPS/SBAS receivers (A314 and A315) are installed at IITRAM provided by SAC ISRO, Ahmedabad for the field trial of the IRNSS. Figure 3 shows antenna set up at terrace of IITRAM with Accord made triband antenna and proposed S Band antenna. Proposed S band antenna is connected to A315,



**Fig. 2** Simulated and measured value of return loss (reflection coefficient) S11 in dB for proposed antenna

**Table 2** Comparison of proposed antenna performance parameters with other S band antenna design from the literature

Reference remarks	Size of ant. mm <sup>2</sup>	Patch type and size	BW %	Gain dBi
Nascetti et al. (2015) S band 2450 MHz	96 * 96	Square 57 * 57 mm <sup>2</sup>	Not metioned	7.3
Pachigolla et al. (2018) ISM band 2.4 GHz	50 * 50	Reactangle 29 * 38 mm <sup>2</sup>	3.75	1.75 for FR4, 4.1 for Arlon
Desai et al. (2018) Transparent Ant MIMO Band 1: 2.4 GHz WLAN Band 2: 3.7 GHz WiMAX	50 * 50	Slotted interconnected ring resonator 24 mm diameter	Single element Band 1: 18.70 Band 2: 21.28 2*1 element Band 1: 11.29 Band 2: 11.64	Single element Band 1: 1.12 Band 2: 2.28 2 * 1 element Band 1: 1.98 Band 2: 2.95
Desai and Upadhyaya (2018) Transpar-ent Ant for smart devices Band 1: 2.4 GHz Band 2: 5.5 GHz	35 * 35	Two over lapping rings 16.2 mm diameter	Band 1: 5.61 Band 2 : 3.62	Band 1: 0.70 Band 2: 1.67
Hussein et al. (2019) S band 2 to 4 GHz	70 * 70	Gear shaped radiating patch 31 mm diameter	2.39	4.27
Proposed antenna	80 * 60	Rectangle 38 * 29 mm <sup>2</sup>	2.41	7.43

**Fig. 3** Antenna set up at terrace of IITRAM

Accord made triband antenna is connected with A314 and GPS antenna is connected to its GPS receiver to collect the positioning data. The antenna is mounted at approximately 25-m height from the ground at IITRAM. The value of

carrier to noise ratio ( $C/N_0$ ) and positioning data have been analyzed to find the signal strength and position accuracy of IRNSS at stationary condition. Table 3 shows the specifications of our experimental setup using which the performance of IRNSS is analyzed with proposed S band antenna and Fig. 4 shows the front view of IRNSS/GPS/SBAS receiver and the computer monitor which shows the value of carrier to noise ratio. The vertical bars shown on both the screen are the amplitude of carrier to noise ratio. One bar indicates one channel received by receiver. As per Fig. 4, six channels are received, and according to atmospheric conditions its amplitude varies frequently.

## Results and discussion

The performance analysis of IRNSS is carried out with three different antennas to evaluate the signal strength in terms of carrier to noise ratio and positioning error. The Accord made triband antenna is connected to A314 IRNSS/GPS/SBAS receiver, the proposed S band antenna is connected to A315 IRNSS/GPS/SBAS receiver to

**Table 3** Specification of experimental setup

System	Indian Regional Navigation Satellite System
Antenna type	Proposed antenna, accord made triband antenna, GPS antenna
Receiver type	A314, A315 accord IRNSS/GPS/SBAS receiver, GPS receiver
Data type	L1, L5, S band
Location	Institute of Infrastructure Technology Research and Management, Ahmedabad, Gujarat, India
Time period	25–31 March 2019



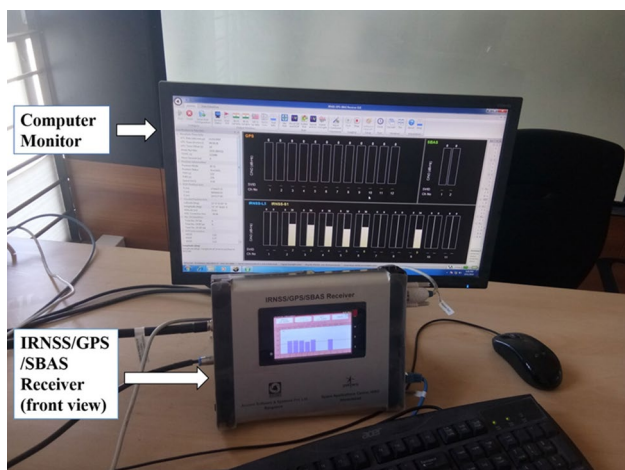


Fig. 4 Experimental setup at lab of IITRAM

Table 4 Carrier to noise ratio for IRNSS channels received by Prototype antenna

Channel no.	PRN	IRNSS+GPS mode (L1+L5+S) (Accord system antenna)		IRNSS S band mode (Prototype antenna)	
		C/N <sub>0</sub> <sup>a</sup> (dB-Hz)	Iono type	C/N <sub>0</sub> <sup>a</sup> (dB-Hz)	Iono type
1	1	38.5	Dual	40.3	Grid
2	2	43.6	Dual	38.9	Grid
3	3	43.0	Dual	37.9	Grid
4	4	42.2	Dual	38.8	Grid
5	5	41.6	Dual	39.6	Grid
6	6	40.3	Dual	41.2	Grid
7	7	39.4	Dual	39.3	Grid

<sup>a</sup>These data are changed continuously with minor variations according to the time

collect real time navigation signal. Table 4 shows the values of carrier to noise ratio in dB and ionospheric delay correction model for Accord made triband antenna and proposed S band antenna. The acceptable value of C/N<sub>0</sub> for IRNSS is greater than 32 dB (Parmar et al. 2015). We found the carrier to noise ratio is between 37 and 41 dB for each channel of satellite and is in acceptable range for

Table 5 Positioning information of IIT RAM with three different antennas

S. no.	Antenna type	Latitude	Longitude	Altitude (m)
1	GPS antenna	23.00440922470994	72.62166043201685	80.69
2	IRNSS dual band + GPS antenna	23.0044004129077	72.6216892127453	78.21
3	IRNSS S band antenna	23.0044471749422	72.6216352279293	79.23

proposed S band prototype antenna. The first and eighth channels are not used for tracking. Channel number two to seven must be tracked and ninth channel can be tracked by adding PRN message comment in IRNSS (ISRO 2019). The positioning information with Accord made triband antenna is 23.0044004129077 and 72.6216892127453 and positioning information with proposed S band antenna is 23.0044471749422 and 72.6216352279293. This shows that the positioning data acquired by S band antenna is in good approximation with the triband antenna. Here, the location coordinates of IITRAM which are 23.00439687 and 72.62182262 have been considered as a golden reference to determine the error in positioning (Rawat et al. 2018). So, the positioning error for S band prototype antenna is found which is less than 10 m. Table 5 shows the latitude, longitude and altitude of IITRAM at stationary point using three different antennas with three different receivers. Table 6 shows that positioning error is 1.9 m for GPS antenna, 1.27 m for Accord made triband antenna (IRNSS dual band + GPS antenna) and 1.29 m for proposed S band antenna. The altitude of the rooftop of the IITRAM is 84.50 m at which the IRNSS antennas are mounted. The altitude error is 3.81 m for GPS antenna, 6.29 m for Accord made triband antenna (IRNSS dual band + GPS antenna) and 5.27 m for proposed S band antenna. For IITRAM, the X position is 1,754,426.3 m, Y position is 5,605,814.42 m and Z position is 2,477,176.59 m, which are considered as a golden reference or reference position stored in IRNSS/GPS/SBAS receiver to calculate the positioning error. The position error is calculated by Eq. (3) and it is verified with the receiver data also. The positioning error can be calculated by the following equation,

$$\text{Error} = \sqrt{(X - X_1)^2 + (Y - Y_1)^2 + (Z - Z_1)^2} \tag{3}$$

where X<sub>1</sub>, Y<sub>1</sub> and Z<sub>1</sub> are the instantaneous value of real time positioning data collected from the IRNSS receiver. The positioning data from receiver is subtracted from Golden reference, X, Y, Z individually, by adding and taking square of its answer and finally finding square root of it, we get the final positioning error at the location of IITRAM for stationary point.

**Table 6** Position and altitude error at IIT RAM with three different antennas

S. no.	Antenna type	Position error (m)	Altitude error (m)
1	GPS antenna	1.9	3.81
2	IRNSS dual band + GPS antenna	1.27	6.29
3	IRNSS S band antenna	1.29	5.27

## Conclusion

In this paper, a compact and low cost circularly polarized Microstrip patch antenna design is proposed for S band application of IRNSS receiver. Initially, the hardware prototype of the proposed antenna is tested using Keysight Field fox Microwave Analyzer N9916A and found to be shifted from the desired resonant frequency of IRNSS receiver. After proper single stub matching, we achieved the desired resonant frequency and bandwidth in the hardware prototype. The positioning information of IITRAM is found by connecting Accord made triband antenna, proposed S band antenna and GPS antenna to individual receivers and compared the error in positioning data. Positioning error is least with Accord triband antenna which receives L1, L5 and S band data and altitude error is least with GPS antenna. A single frequency Grid model for ionosphere correction is applied for S band operation to correct the ionosphere delay and provides precise positioning data for single frequency user. Results reveal that the positioning error of 1.29 m for the single frequency operation using proposed prototype S band antenna is achieved, which are less than 10 m and fulfill the objective of Indian Regional Navigation Satellite System. Further, the signal strength, carrier to noise ratio and positioning data are found with good accuracy with proposed single band (S band) antenna also. This compact antenna can be used to find PVT and proved a good choice in the low cost for single band navigation solutions.

**Acknowledgements** We would like to express thanks to SAC ISRO (Indian Space Research Organization), Ahmedabad to provide IRNSS/GPS/SBAS receiver for the field trails. We are also thankful to the team of Accord Systems, Bengaluru for providing technical support.

## Compliance with ethical standards

**Conflict of interest** The authors hereby confirm that there is no conflict of interest with respect to the current manuscript.

## References

Bilotti F, Vegni C (2010) Design of high-performing microstrip receiving GPS antennas with multiple feeds. *IEEE Antennas Wirel Propag Lett* 9:248–251

- Desai A, Upadhyaya T (2018) Transparent dual band antenna with negative material loading for smart devices. *Microw Opt Technol Lett* 60(11):2805–2811. <https://doi.org/10.1002/mop.31474>
- Desai A, Upadhyaya T, Palandoken M, Gocen C (2018) Dual band transparent antenna for wireless MIMO system applications. *Microw Opt Technol Lett* 61(7):1845–1856
- Hussein AH, Abdullah HH, Attia MA, Abada AM (2019) S-band compact microstrip full-duplex tx/rx patch antenna with high isolation. *IEEE Antennas Wirel Propag Lett* 18(10):2090–2094. <https://doi.org/10.1109/LAWP.2019.2937769>
- ISRO (2019) Indian Regional Navigation Satellite System. [https://www.isro.gov.in/sites/default/files/irnss\\_pdf](https://www.isro.gov.in/sites/default/files/irnss_pdf). Accessed 07/02/2019
- Kiran B, Raghu N, Manjunatha KN (2017) A comparative study and performance analysis using IRNSS and hybrid satellites. In: Sathapathy S, Prasad V, Rani B, Udgata S, Raju K (eds) *Proceedings of the first international conference on computational intelligence and informatics. Advances in Intelligent Systems and Computing*, vol 507. Springer, Singapore
- Kumar G, Ray KR (2003) *Broadband microstrip antennas*. Artech House, Boston
- Nascetti A, Pittella E, Teofilatto P, Pisa S (2015) High-gain s-band patch antenna system for earth-observation CubeSat satellites. *IEEE Antennas Wirel Propag Lett* 14:434–437
- Nasimuddin Chen ZN, Qing X (2012) A compact circularly polarized crossshaped slotted microstrip antenna. *IEEE Trans Antennas Propag* 60(3):1584–1588. <https://doi.org/10.1109/TAP.2011.2180334>
- Pachigolla SSY, Dab V, Chatterjee A, Kundu S (2018) A compact rectangular microstrip patch antenna for 2.4 GHz ISM band applications. In: 2018 IEEE Indian conference on antennas and propagation (InCAP), pp 1–3
- Parmar S, Dalal U, Pathak KN (2015) Detecting ionospheric irregularities using empirical mode decomposition of TEC for IRNSS signals, at SVNIT, Surat, India. *Commun Appl Electron: CAE* 7(8):22–29
- Rawat A, Savaliya J, Chhabhaya D (2018) Field trial of IRNSS receiver. *Microw Opt Technol Lett* 61(5):1149–1153. <https://doi.org/10.1002/mop.31746>
- Ray D, Kumar G (2000) Tuneable and dual-band circular microstrip antenna with stubs. *IEEE Trans Antennas Propag* 48:1036–1039. <https://doi.org/10.1109/8.876321>
- Reddy BS, Kumar VS, Srinivasan VV, Mehta Y (2015) Dual band circularly polarized microstrip antenna for IRNSS reference receiver. In: 2015 IEEE MTT-S international microwave and RF conference (IMaRC), pp 279–282. <https://doi.org/10.1109/imarc.2015.7411432>
- Roy SK, Jha L (2019) Effects of tuning stub on microstrip patch antenna. <https://www.niscair.res.in/sciencecommunication/2005>
- Sahal M, Tiwari V (2016) Review of circular polarization techniques for design of microstrip patch antenna, pp 663–669. <https://doi.org/10.1007/978-81-322-2638-374>
- Sharma PC, Gupta KC (1983) Analysis and optimized design of single feed circularly polarized microstrip antennas. *IEEE Trans Antennas Propag* 31:949–955. <https://doi.org/10.1109/tap.1983.1143162>
- Wu S, Liu S, Guo Z (2010) Coaxial probe-fed circularly polarized microstrip antenna for Beidou RDSS applications. In: 2010 international conference on microwave and millimeter wave technology, pp 297–299

# Acta Geophysica

Volume 68  
Number 4  
2020

## RESEARCH ARTICLES - SOLID EARTH SCIENCES

### **Seismicity analysis of selected faults in Makran Southern Pakistan**

M.J. Khan · M. Ali · M. Xu · M. Khan 965

### **The use of QLARM to estimate seismic risk in Kirghizstan at the regional and city scales**

P. Rosset · S. Tolis · M. Wyss 979

### **Numerical modelling of the near-field velocity pulse-like ground motions of the Northridge earthquake**

Q. Luo · F. Dai · Y. Liu · M. Gao 993

## RESEARCH ARTICLES - APPLIED GEOPHYSICS

### **Thin interbed AVA inversion based on a fast algorithm for reflectivity**

Z. Yang · J. Lu 1007

### **Three-dimensional angle-domain double-square-root migration in VTI media for the large-scale wide-azimuth seismic data**

C. Wu · B. Feng · H. Wang · T. Wang 1021

### **Stable absorption compensation with lateral constraint**

X. Ma · G. Li · H. Li · J. Li · X. Fan 1039

### **Three-dimensional magnetotelluric inversion using L-BFGS**

L. Lu · K. Wang · H. Tan · Q. Li 1049

### **Prestack AVO inversion for brittleness index of shale based on BI\_Zoeppritz equation and NSGA II**

C. Bi · Y. Wang · W. Xie · W. Sun · W. Liu 1067

### **Interpretation of gravity anomaly over 2D vertical and horizontal thin sheet with finite length and width**

A. Biswas 1083

## REVIEW ARTICLE - ANTHROPOGENIC HAZARD

### **Koyna earthquakes: a review of the mechanisms of reservoir-triggered seismicity and slip tendency analysis of subsurface faults**

D. Das · J. Mallik 1097

## RESEARCH ARTICLES - HYDROLOGY

### **Modelling reference evapotranspiration by combining neuro-fuzzy and evolutionary strategies**

M. Alizamir · O. Kisi · R. Muhammad Adnan ·

A. Kuriqi 1113

### **Mapping shoreline change using machine learning: a case study from the eastern Indian coast**

L. Kumar · M.S. Afzal · M.M. Afzal 1127

### **Dynamics of thin disk settling in two-layered fluid with density transition**

M.M. Mrokowska 1145

(Contents continued on inside back cover)