

ISSN 2300-5149

PRZEGLĄD TELEINFORMATYCZNY

W. Kwiatkowski Rekomendacje jako wynik oceny preferencji na podstawie wskazanych przykładów	3
M. Olchowik Analysis and classification of chosen social engineering methods in cybersecurity	17
A. Arciuch Ocena dokładności i powtarzalności pomiaru parametru RSSI w systemach BLE.....	31
A.M. Donigiewicz Recenzja książki „Komunikacja sieciowa: źródła informacji big data” autorstwa Dariusza Jarugi	41
Information for Authors – rules of papers preparation and reviewing for Teleinformatics Review	49
Informacje dla autorów – zasady przygotowania tekstu i recenzowania artykułów do Przeglądu Teleinformatycznego	51

PRZEGLĄD TELEINFORMATYCZNY
TELEINFORMATICS REVIEW

Dawniej: BIULETYN INSTYTUTU AUTOMATYKI I ROBOTYKI WAT
(ISSN 1427-3578)
Ukazuje się od 1995 r.

RADA NAUKOWA

Lt. Col. Janos Balogh MSc
dr hab. inż. Antoni M. Donigiewicz – przewodniczący
prof. Hacene Fouchal, PhD
prof. Lech J. Janczewski, DEng
prof. dr hab. inż. Włodzimierz Kwiatkowski
prof. dr hab. inż. Bohdan Macukow
Lt. Col. Lajos Mucha PhD
prof. ing. Vladimír Olej, CSc.

ADRES REDAKCJI

Redakcja Przeglądu Teleinformatycznego
00-908 Warszawa, ul. gen. Sylwestra Kaliskiego 2
tel. 261 83 87 03, fax. 261 83 71 44
e-mail: pt [at] ita.wat.edu.pl

WWW: <https://przegladteleinformatyczny.publisherspanel.com/>
<http://przeglad.ita.wat.edu.pl/>

Wersją pierwotną czasopisma jest wersja elektroniczna

REDAKTOR NACZELNY:

Antoni Donigiewicz

REDAKTOR WYDANIA

Antoni Donigiewicz

OPRACOWANIE STYLISTYCZNE

Renata Borkowska

PROJEKT OKŁADKI

Barbara Chruszczyk

WYDAWCA: Instytut Teleinformatyki i Cyberbezpieczeństwa WAT

ISSN 2300-5149
ISSN 2353-9836 (on-line)

Rekomendacje jako wynik oceny preferencji na podstawie wskazanych przykładów

Włodzimierz KWIATKOWSKI

Instytut Teleinformatyki i Cyberbezpieczeństwa, Wydział Cybernetyki, WAT,
ul. gen. S. Kaliskiego 2, 00-908 Warszawa
wlodzimierz.kwiatkowski@wat.edu.pl

STRESZCZENIE: Rozpatrywany jest problem wyznaczania rekomendacji na podstawie wskazanych przykładów decyzji akceptowalnych i przykładów decyzji nieakceptowalnych. Wskazanie przez decydenta tych przykładów jest podstawą oceny jego preferencji. Istota przedstawionego rozwiązania polega na określeniu preferencji jako klastra wyznaczonego poprzez uzupełnianie wskazanych przykładów. W artykule zaproponowano procedurę kolejnych przybliżeń opartą na rozwiązaniach zadania klasyfikacji na podstawie zadanych przykładów.

SŁOWA KLUCZOWE: rekomendacja, preferencje, eksploracja danych, klasyfikacja, grupowanie

1. Wprowadzenie

Rozpatrywane zadanie wyznaczenia rekomendacji polega na sprecyzowaniu, które decyzje z danego zbioru są zgodne z preferencjami decydenta. Preferencje decydenta są wyrażane jako zbiór decyzji akceptowalnych przez niego. Sposobem poznania preferencji decydenta jest analiza wskazywanych przez niego przykładów. Wyznaczony na tej podstawie zbiór decyzji akceptowalnych jest traktowany jako rekomendacja dla decydenta.

Definiowanie cech i wskaźników jakości decyzji w sposób niezależny od konkretnego aktu wyboru decydenta praktycznie oznacza, że decydent ma preferencje arbitralnie narzucone. Wnioskowanie o preferencjach decydenta tylko na podstawie wskazywanych przez niego przykładów decyzji ocenionych pozytywnie (akceptowalnych) wyklucza taką sytuację. Takie wnioskowanie oznacza także, że ewaluacja decyzji dokonywana jest bezpośrednio na podstawie

ich charakterystyk (np. pomiarów, obserwacji), a nie na podstawie narzuconych cech i wskaźników jakościowych.

Sformułowanie problemu wyznaczania rekomendacji na podstawie wskazanych przez decydenta przykładów decyzji jest przedstawione w pracy [3]. Przyjęta tam metoda ewaluacji opiera się na wyznaczaniu w przestrzeni cech odległości analizowanych decyzji od decyzji wskazanych.

Omawiane w niniejszym artykule zadanie różni się od przedstawionego w publikacji [3]. Różnica wynika z innej interpretacji wskazywanych przykładów. W pracy [3] wskazane przez decydenta przykłady są traktowane jako deklaracja jego preferencji. W konsekwencji tego podane przykłady stanowią ustalony zbiór wzorców. W niniejszym artykule wskazywane decyzje są interpretowane jako deklaracja niepełna. Istota przedstawianego tu rozwiązania polega na określeniu preferencji decydenta jako klastra (skupienia) wyznaczonego poprzez uzupełnianie wskazanych przykładów.

W przypadku istnienia sprzężenia zwrotnego wzbogacanie zbioru przykładów polega na cyklicznym zapoznawaniu decydenta z proponowaną rekomendacją i uzyskiwaniu od niego dodatkowych wskazówek. Istota metody proponowanej w niniejszym artykule polega na wykorzystaniu idei stopniowego wzbogacania zbioru decyzji wskazanych przez decydenta. Przyjmuje się przy tym założenie, że proces modyfikacji powinien następować bez udziału decydenta. Podstawą możliwości realizacji tej idei jest analiza wskazywanych przykładów na tle wszystkich rozpatrywanych decyzji. Taką analizę umożliwia przedstawiona w pracy [4] metoda regularyzacji zadań klasyfikacji.

Rozpatrywany w artykule problem jest sformułowany jako poszukiwanie metody wyznaczania rekomendacji na podstawie wskazania przez decydenta przykładów decyzji akceptowalnych, a także przykładów decyzji nieakceptowalnych. Celem wskazywania przykładów decyzji nieakceptowalnych jest racjonalne ograniczanie liczebności rekomendowanych decyzji.

2. Prace związane

Sformułowane zadanie wyznaczania rekomendacji można zaliczyć do projektowania systemów eksploracji danych, których wyróżnikiem jest wyszukiwanie informacji według zgłoszonego zapotrzebowania użytkownika. Przyjęte założenia powodują, że poszukiwane jest rozwiązanie polegające na filtrowaniu decyzji opartym zarówno na treści (ang. *content based filtering*), jak i na współpracy (ang. *collaborative filtering*) [6].

Podstawowym przykładem iteracyjnego grupowania poprzez kolejne modyfikacje wyników jest algorytm ISODATA (ang. *Iterative Self-Organizing Data Analysis Techniques*) [1]. Określenie *algorytm ISODATA* jest najczęściej

rozumiane jako szczególna metoda nienadzorowanej klasyfikacji [2]. Obecnie algorytm ISODATA jest wykorzystywany np. przy klasyfikacji obrazów multi-spektralnych¹.

Zasadniczy kłopot pojawiający się przy klasyfikacji na podstawie wskazywanych przykładów wynika z faktu, że wskazywane przykłady generują podprzestrzeń, której wymiar jest mniejszy od wymiaru przestrzeni cech. Problem wynikający z faktu, że liczba wskazanych przykładów wzorców jest mała względem liczby współrzędnych wektora cech, jest rozpatrywany w pracy [5]. Zaproponowane są tam dwie metody optymalizacji oparte na wyznaczaniu rzutów wektorów cech na podprzestrzeń wzorców. Wyróżnikiem pierwszej metody jest wykorzystywanie odległości wektora cech od podprzestrzeni wzorców. Druga metoda polega na przeniesieniu zadania optymalizacji do podprzestrzeni wzorców. Przedstawiona w publikacji [4] metoda regularyzacji jest rozwiązaniem kompromisowym i pozwala efektywnie wykorzystywać wskazania decydenta także w przypadkach osobliwych macierzy kowariancji cech wskazanych przykładów.

3. Klasyfikacja na podstawie zadanych wzorców klas

Dany jest zbiór decyzji ponumerowany od 1 do N . Dla każdej decyzji znany jest jej wektor cech. Dla decyzji o numerze k stosować będziemy następujące oznaczenie wektora cech:

$$\mathbf{a}_k = [a_{1,k}, a_{2,k}, \dots, a_{L,k}]^T, \quad \mathbf{a}_k \in R^L \quad (1)$$

Każda współrzędna $a_{l,k}$ jest liczbą rzeczywistą, a parametr L określa liczbę współrzędnych wektora cech. Wektory cech zadanego zbioru decyzji zestawiamy w postaci następującej macierzy:

$$\mathbf{A} = [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_N], \quad \mathbf{a}_k \in R^L \quad (2)$$

Macierz kowariancji wektorów cech wyznaczana jest następująco:

$$\mathbf{R} = \frac{1}{N-1} \sum_{k=1}^N (\mathbf{a}_k - \bar{\mathbf{a}})(\mathbf{a}_k - \bar{\mathbf{a}})^T \quad (3)$$

gdzie:

$$\bar{\mathbf{a}} = \frac{1}{N} \sum_{k=1}^N \mathbf{a}_k \quad (4)$$

¹ Na podstawie: https://en.wikipedia.org/wiki/Multispectral_pattern_recognition. Dostęp: 23.08.2021.

Przyjmiemy dalej, że macierz kowariancji wektorów cech jest nieosobliwa:

$$\det(\mathbf{R}) \neq 0 \quad (5)$$

Odległość pomiędzy wektorami \mathbf{x} , \mathbf{y} przestrzeni cech R^L będziemy wyznaczać w sposób uwzględniający wielkość rozrzutu (rozproszenia) wartości współrzędnych oraz ich wzajemną korelację. Wymagania te spełnia odległość Mahalanobisa, która jest określona wzorem:

$$d_e(\mathbf{x}, \mathbf{y}) = \sqrt{(\mathbf{x} - \mathbf{y})^T \mathbf{R}^{-1} (\mathbf{x} - \mathbf{y})}, \quad \mathbf{x}, \mathbf{y} \in R^L \quad (6)$$

Wskazania przykładów określających wzorzec klasy o indeksie $h \in \{1, 2, \dots, H\}$ (gdzie: H – liczba klas) będziemy dokonywać przez podanie odpowiedniego zbioru indeksów W_h . Liczbę elementów wzorca o indeksie h oznaczamy jako

$$N_h = \|W_h\| \quad (7)$$

Wzorzec klasy o indeksie h jest reprezentowany przez następujący zbiór punktów (klastrów) w przestrzeni cech:

$$C(W_h) = \{\mathbf{a}_k \in R^L : k \in W_h\} \quad (8)$$

Wnioskowanie o podobieństwie cechy \mathbf{x} do wzorca o indeksie h opiera się na określeniu odległości $D_e(\mathbf{x}, C(W_h))$ punktu \mathbf{x} od klastra $C(W_h)$. Przykładowo, wybierając metodę centroidalną wyznaczania odległości między klastrami, otrzymujemy zależność:

$$D_e(\mathbf{x}, C(W_h)) = d_e(\mathbf{x}, \bar{\mathbf{w}}_h) = \sqrt{(\mathbf{x} - \bar{\mathbf{w}}_h)^T \mathbf{R}^{-1} (\mathbf{x} - \bar{\mathbf{w}}_h)} \quad (9)$$

gdzie:

$$\bar{\mathbf{w}}_h = \frac{1}{N_h} \sum_{j \in W_h} \mathbf{a}_j \quad (10)$$

Klasyfikacja oparta na wykorzystywaniu metryki (6) nazywana jest środowiskową [4].

Stosowanie klasyfikacji środowiskowej znajduje uzasadnienie wtedy, gdy cechy wszystkich wzorców są jednorodnie w następującym sensie: odpowiednie klastry różnią się wartościami oczekiwanymi, a odpowiadające im macierze kowariancji są jednakowe. W przypadku, gdy macierze kowariancji wzorców różnią się, zalecane jest zróżnicowanie sposobu pomiaru odległości stosownie do macierzy kowariancji poszczególnych wzorców [2].

Macierz kowariancji wyznaczoną na podstawie przykładów wzorca o indeksie h oznaczmy następująco:

$$\mathbf{R}_h = \frac{1}{N_h - 1} \sum_{j \in W_h} (\mathbf{a}_j - \bar{\mathbf{w}}_h)(\mathbf{a}_j - \bar{\mathbf{w}}_h)^T \quad (11)$$

Odległość pomiędzy wektorami \mathbf{x} , \mathbf{y} przestrzeni cech R^L zadaną wzorem:

$$d_h(\mathbf{x}, \mathbf{y}) = \sqrt{(\mathbf{x} - \mathbf{y})^T \mathbf{R}_h^{-1} (\mathbf{x} - \mathbf{y})}, \quad \mathbf{x}, \mathbf{y} \in R^L \quad (12)$$

nazywa się dopasowaną do wzorca o indeksie h [2]. Podobnie nazywać będziemy odległość między cechą \mathbf{x} a klastrem $C(W_h)$. Przykładowo dla centroidalnej metody grupowania odległość ta jest określona wzorem:

$$D_h(\mathbf{x}, C(W_h)) = d_h(\mathbf{x}, \bar{\mathbf{w}}_h) = \sqrt{(\mathbf{x} - \bar{\mathbf{w}}_h)^T \mathbf{R}_h^{-1} (\mathbf{x} - \bar{\mathbf{w}}_h)} \quad (13)$$

Klasyfikacja względem zadanych wzorców polega na przyporządkowaniu analizowanej decyzji o indeksie k do klasy o indeksie h wtedy, jeśli [2]

$$D_h(\mathbf{a}_k, C(W_h)) = \min_{j=1,2,\dots,H} D_j(\mathbf{a}_k, C(W_j)) \quad (14)$$

Potrzeba regularyzacji występuje, gdy macierze kowariancji cech wzorców są osobliwe lub źle uwarunkowane (występuje duża rozpiętość między ich wartościami własnymi, a ich wyznaczniki są bliskie zeru). Zaproponowana w pracy [4] metoda regularyzacji polega wykorzystaniu następującej metryki:

$$d_{h,\rho}(\mathbf{x}, \mathbf{y}) = \sqrt{(\mathbf{x} - \mathbf{y})^T \mathbf{R}_{h,\rho}^{-1} (\mathbf{x} - \mathbf{y})}, \quad \mathbf{x}, \mathbf{y} \in R^L \quad (15)$$

przy czym:

$$\mathbf{R}_{h,\rho} = (1 - \rho)\mathbf{R}_h + \rho\mathbf{R} \quad (16)$$

gdzie: $\rho \in [0,1]$ – współczynnik regularyzacji. Metryka $d_{h,\rho}(\mathbf{x}, \mathbf{y})$ różni się od metryki $d_h(\mathbf{x}, \mathbf{y})$ zastąpieniem macierzy kowariancji \mathbf{R}_h kombinacją wypukłą tej macierzy i macierzy kowariancji \mathbf{R} . Wartość $\rho=0$ oznacza brak regularyzacji i klasyfikację dopasowaną, wartość $\rho=1$ oznacza przejście do klasyfikacji środowiskowej.

4. Rekomendacje na podstawie klasyfikacji względem wzorców

Podstawowe założenie poczynione przy wyznaczaniu przedstawianego w tym punkcie rozwiązania jest następujące: wskazane przykłady decyzji akceptowalnych są traktowane jako wzorzec decyzji akceptowalnych, a wskazane przykłady decyzji nieakceptowalnych – jako wzorcowe dla odpowiadającej im klasy. Możliwe jest zastosowanie dwóch metod klasyfikacji: środowiskowej i dopasowanej do wzorców [3]. W obu przypadkach rozwiązanie zadania klasyfikacji można bezpośrednio wykorzystać do wyznaczania rekomendacji. Uzyskiwane rozwiązanie przedstawia podział zbioru decyzji na dwie rozdzielne klasy: decyzji akceptowalnych i decyzji nieakceptowalnych. Wyznaczoną klasę

decyzji akceptowalnych można traktować jako rekomendację dla decydenta. Metoda ta obarczona jest wrodzoną wadą: rekomendowane są wszystkie decyzje, którym jest bliżej do przykładów akceptowalnych niż do przykładów nieakceptowalnych. Można tę wadę osłabić, zadając odpowiednio małą liczebność zbioru rekomendowanych decyzji lub nakładając na decyzje rekomendowane jakościowe ograniczenia. Takie ograniczenia liczebności zbioru decyzji rekomendowanych są możliwe na podstawie uszeregowania elementów wyznaczonej klasy na podstawie ich odległości od swojego wzorca.

Jest oczywiste, że rekomendacje można wyznaczać zarówno dla decyzji akceptowalnych, jak i nieakceptowalnych. Jest to niewątpliwa zaleta wykorzystywania zadania klasyfikacji. Podział zbioru decyzji na dwie przeciwstawne klasy jest najbardziej ostrym opisem preferencji decydenta.

5. Rekomendacje na podstawie szeregowania decyzji

Podobnie jak poprzednio, wskazane przykłady decyzji akceptowalnych są traktowane jako wzorce decyzji akceptowalnych. Podstawą wyznaczenia rekomendacji jest uporządkowanie (uszeregowanie) zbioru decyzji według odległości decyzji od klastra wzorców, rozpoczynając od decyzji położonej w odległości najmniejszej. Do zbioru decyzji rekomendowanych kwalifikowane są kolejne decyzje zgodnie z tym uszeregowaniem. Jest oczywiste, że wskazane przykłady decyzji nieakceptowalnych należy wykluczyć.

Bardziej wnikliwe wnioskowanie o wykluczeniu decyzji w procesie kwalifikowania do rekomendacji jest możliwe na podstawie uporządkowania (uszeregowania) zbioru decyzji względem wzorca decyzji nieakceptowalnych. W tym przypadku należy więc wykonać uporządkowanie zbioru wszystkich decyzji zarówno względem wzorców decyzji akceptowalnych, jak i uporządkowanie tego zbioru względem wzorców decyzji nieakceptowalnych. Konfrontacja tych dwóch przeciwstawnych uszeregowień powinna umożliwić racjonalne ograniczanie liczebności zarówno zbioru rekomendowanych decyzji akceptowalnych, jak i zbioru rekomendowanych decyzji nieakceptowalnych. Proponowana idea ograniczania sekwencji jest następująca. Po napotkaniu w szeregu decyzji uporządkowanych według wzorców akceptowalnych decyzji wskazanej jako przykład (wzorzec) decyzji nieakceptowalnej z procesu rekomendowania wyklucza się także wszystkie decyzje następne w szeregu (usytuowane dalej od klastra wzorców decyzji akceptowalnych). Analogicznie można ograniczać rekomendacje decyzji nieakceptowalnych.

Realizacja przedstawionej wyżej idei polega na niezależnym przeprowadzeniu dwóch procesów kwalifikowania decyzji: do zbioru decyzji rekomendowanych jako akceptowalne i do zbioru decyzji rekomendowanych jako

nieakceptowalne. W obu przypadkach proces kwalifikowania do odpowiedniej rekomendacji należy zakończyć po napotkaniu w szeregu decyzji wskazanej jako kontrprzykład. Po wykluczeniu elementów wspólnych uzyskane klastry mogą być traktowane jako odpowiednie, przeciwstawne rekomendacje.

6. Metoda iteracyjna wyznaczania rekomendacji na podstawie szeregowania

Wskazanie wzorcowych przykładów przez decydenta pośrednio daje informację o tym, które współrzędne wektora cech i jakie ich wartości są dla decydenta istotne.

Poszerzenie wymiaru podprzestrzeni generowanej przez wzorce można uzyskać poprzez zwiększenie liczebności zbioru wzorców. Proponujemy w tym celu, aby rekomendacje uzyskane na podstawie wskazanych przez decydenta przykładów zinterpretować jako nowy zbiór decyzji definiujących odpowiednią klasę decyzji: akceptowalnych bądź nieakceptowalnych. Uznawanie rekomendacji za nowy wzorec klasy można powtarzać.

Ideę proponowanej metody można przedstawić jako iteracyjne uzupełnianie klastra wzorców decyzji akceptowalnych i klastra wzorców decyzji nieakceptowalnych. Uzupełnianie polega na wyznaczeniu w każdej iteracji klastra decyzji rekomendowanych jako akceptowalne i klastra decyzji rekomendowanych jako nieakceptowalne. W kolejnej iteracji wyznaczone klastry decyzji rekomendowanych stają się odpowiednimi klastrami wzorców. Efektywną metodą wyznaczania przeciwnych rekomendacji jest opisana wcześniej metoda oparta na podwójnym szeregowaniu decyzji: względem wzorców decyzji akceptowalnych oraz względem wzorców decyzji nieakceptowalnych. Istotna dla procesu rekomendowania jest eliminacja elementów wspólnych w obu sekwencjach. Procedurę tę można uzupełnić ograniczeniami liczebności zbiorów decyzji rekomendowanych oraz ograniczeniami natury jakościowej. Efektywnym działaniem zapewniającym te efekty jest ograniczanie maksymalnej odległości analizowanej decyzji od klastra wzorcowego.

Oczekiwanym rezultatem opisanego iteracyjnego procesu jest uzyskanie ustalonych zbiorów rekomendowanych decyzji. Zbiory te nie powinny ulegać zmianie w kolejnych iteracjach. Uzyskane, ustalone klastry (skupienia) stanowią bezpośredni opis preferencji decydenta, zarówno akceptowalności, jak i nieakceptowalności decyzji. Bezpośredni opis preferencji oznacza tu wyliczenie (enumerację) wszystkich preferowanych decyzji. Utożsamienie tak rozumianego opisu preferencji z odpowiednią rekomendacją jest naturalną konsekwencją.

7. Eksperyment obliczeniowy

7.1. Przedmiot i cel badań

Celem badań było eksperymentalne potwierdzenie uzyskiwania użytecznych rekomendacji poprzez iteracyjne uzupełniania wskazanych przez decydenta przeciwstawnych przykładów. Istotnym badanym problemem była ocena procesu uzyskiwania rekomendacji, a w szczególności potwierdzenie uzyskiwania rekomendacji ustalonych (tj. nie zmieniających się w kolejnych iteracjach).

Użyte w badaniach dane pomiarowe zostały pobrane z archiwalnych baz danych stacji IMGW². W wybranym do analizy zbiorze pomiarów „decyzja” oznacza pojedynczą stację i jest scharakteryzowana wektorem $L = 2$ pomiarów. Analizowana baza danych zawiera wyniki pomiarów 62 stacji. Pomiarzy zostały wykonane w tym samym dniu.

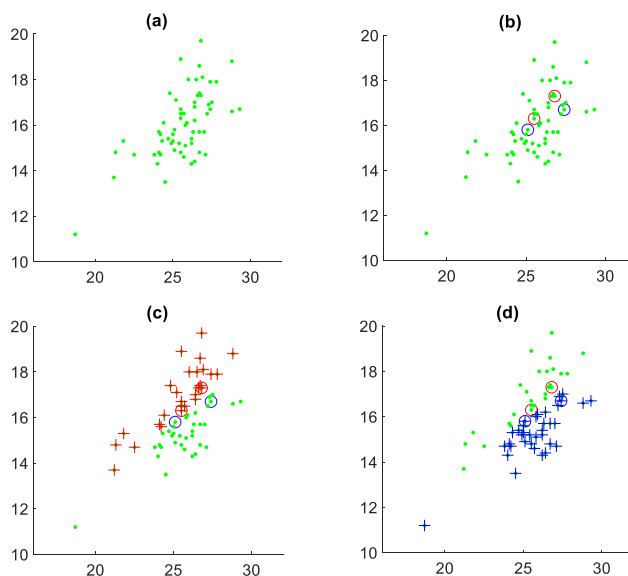
Przyjęty do obliczeń wymiar przestrzeni cech zapewnia czytelną wizualizację wyników. Szczególnie ilustracja przypadku małej liczby wskazanych przykładów powinna być łatwo interpretowalna (wskazanie tylko dwóch przykładów sugeruje ograniczenie preferencji do punktów leżących na prostej przechodzącej przez wektory cech wskazanych przykładów).

7.2. Wyniki wyznaczania rekomendacji środowiskowych

Wyniki wyznaczania rekomendacji na podstawie klasyfikacji środowiskowej stanowią punkt wyjścia do oceny proponowanej w artykule metody.

Źródłem danych jest macierz złożona z 62 wektorów pomiarów wykonanych przez poszczególne stacje. Wektory te określają środowisko eksperymentu. Dla przyjętych do obliczeń danych macierz kowariancji wektorów pomiarowych jest nieosobliwa. Umożliwia to oparcie obliczeń na metryce zdefiniowanej wzorem (6).

² Źródło danych: *IMGW Dane pomiarowo-obszaryjne. Dane meteorologiczne, dobowe, klimatyczne*. Plik: k_d_07_2020.csv (2020_07_k.zip). URL: https://dane.imgw.pl/data/dane_pomiarowo_obszaryjne/dane_meteorologiczne/dobowe/klimat/2020/. Dostęp: 23.08.2021.

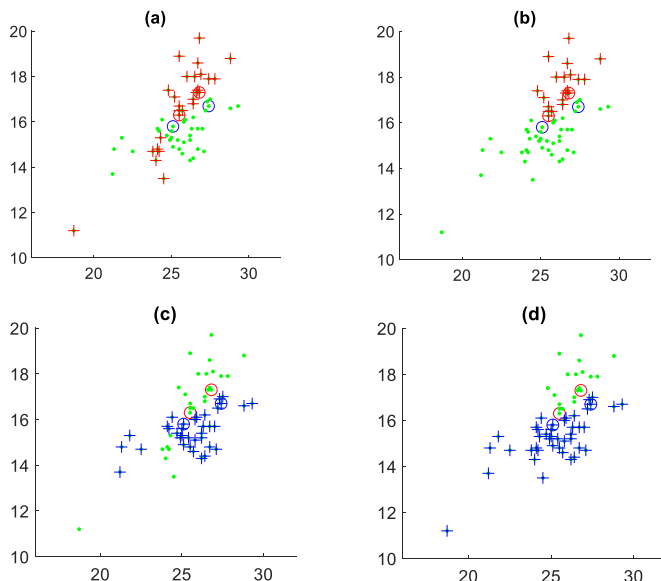


Rys. 1. Wizualizacja wyznaczania rekomendacji środowiskowych. (a) Decyzje środowiska w przestrzeni cech. (b) Wskazane przykłady (wzorce). (c) Rekomendacje dla klasy A (decyzji akceptowalnych). (d) Rekomendacje dla klasy NA (decyzji nieakceptowalnych). Punkty zielone oznaczają wektory cech decyzji środowiska. Wektory cech wskazanych (wzorcowych) decyzji oznaczono kółkami: dla klasy A kolorem czerwonym, dla klasy NA kolorem niebieskim. Wektory cech decyzji rekomendowanych oznaczono znakiem plus w odpowiednim kolorze

Klasyfikacja została przeprowadzona na podstawie wskazanych przykładów: klasy A (decyzji akceptowalnych) oraz wskazanych przykładów klasy NA (decyzji nieakceptowalnych). Przyjęto, że liczba zarówno przykładów klasy A, jak i przykładów klasy NA jest równa 2. Założenie to uniemożliwia bezpośrednie wykorzystywanie odległości dopasowanych do poszczególnych klas (macierze kowariancji dla wskazywanych przykładów klas są osobliwe). Można zauważyć, że w ogólnym przypadku przestrzeni cech R^L do wykonania klasyfikacji dopasowanej potrzebne jest wskazanie co najmniej $L + 1$ przykładów (w dwuwymiarowej przestrzeni cech pożądane jest wskazanie co najmniej trzech przykładów). W praktyce dla dużych wartości L często okazuje się to istotnym problemem, zwłaszcza w przypadku uzyskiwania przykładów w wyniku współpracy z decydującym-człowiekiem.

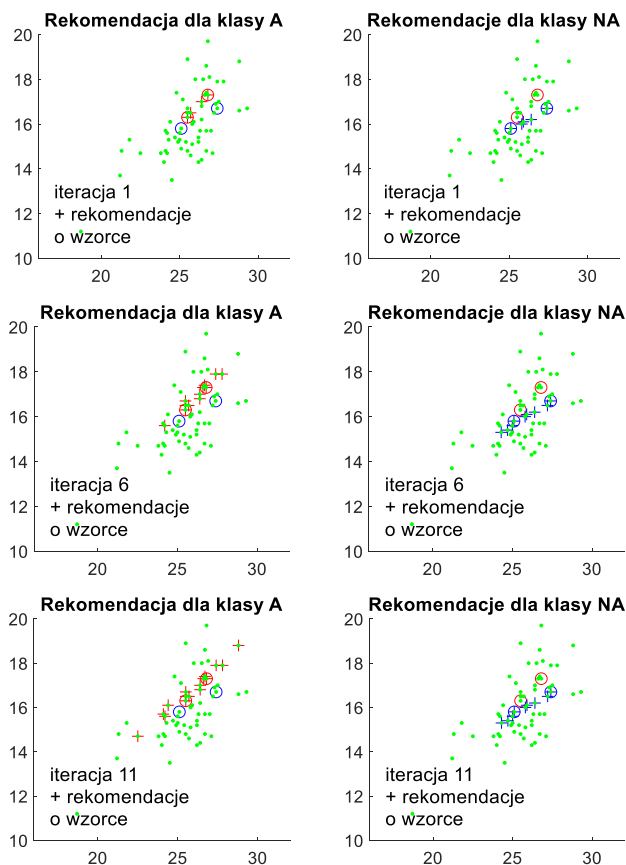
Obliczenia wykonano dla klasyfikacji środowiskowej oraz – po zastosowaniu regularyzacji – dla klasyfikacji dopasowanej do wzorców. Wyniki klasyfikacji środowiskowej są przedstawione na rys. 1. Przyjęte do obliczeń przykłady wzorców zostały specjalnie dobrane na potrzeby wizualizacji. Przedstawione zobrazowanie ukazuje przypadek, kiedy każda analizowana

decyzja jest kwalifikowana jako rekomendacja decyzji akceptowalnej bądź nieakceptowalnej. Mimo oczywistego podejrzenia, że decydent, podając swoje nieliczne przykłady, raczej nie określił tak szeroko swoich preferencji.



Rys. 2. Wizualizacja wyznaczania rekomendacji dopasowanych na podstawie zadania regularyzowanego. (a) Rekomendacje dla klasy A wyznaczone dla współczynnika regularyzacji $\rho=10^{-10}$. (b) Rekomendacje dla klasy A wyznaczone dla współczynnika regularyzacji $\rho=10^{-1}$. (c) Rekomendacje dla klasy NA wyznaczone dla współczynnika regularyzacji $\rho=10^{-10}$. (d) Rekomendacje dla klasy NA wyznaczone dla współczynnika regularyzacji $\rho=10^{-1}$. Punkty zielone oznaczają wektory cech decyzji środowiska. Wektory cech wskazanych (wzorcowych) decyzji oznaczono kółkami: dla klasy A kolorem czerwonym, dla klasy NA kolorem niebieskim. Wektory cech decyzji rekomendowanych oznaczono znakiem plus w odpowiednim kolorze

Na rys. 2 przedstawiono wyniki obliczeń dla klasyfikacji dopasowanej do wzorców. Przyjęta do obliczeń liczba wskazywanych przykładów (dwa dla każdej klasy) powoduje konieczność wykorzystywania regularyzacji zadania. Wskazane wektory cech decyzji akceptowalnych i ich kontrprzykłady są położone blisko siebie, a proste generowane przez odpowiednie punkty przecinają się, tak aby pokazać wpływ współczynnika regularyzacji na uzyskiwane wyniki. Także w tym przypadku każda analizowana decyzja jest kwalifikowana jako rekomendacja decyzji akceptowalnej bądź nieakceptowalnej.



Rys. 3. Wizualizacja wyznaczania rekomendacji metodą iteracyjną na podstawie regularyzowanego zadania klasyfikacji. Punkty zielone oznaczają wektory cech decyzyjnego środowiska. Wektory cech wskazanych (wzorcowych) decyzji oznaczono kółkami: dla klasy A kolorem czerwonym, dla klasy NA kolorem niebieskim. Wektory cech decyzji rekomendowanych oznaczono znakiem plus w odpowiednim kolorze

7.3. Wyniki iteracyjnego wyznaczania rekomendacji metodą iteracyjnego uzupełniania wskazanych przykładów

W odróżnieniu od poprzednio przedstawianych obliczeń, wskazania decydenta nie są traktowane jako zamknięte listy wzorców klas, a jedynie ich przykłady. Podstawą zakwalifikowania do decyzji rekomendowanych jest szeregowanie decyzji względem wskazanych przykładów. Proces kwalifikowania kolejnych decyzji jest kończony w momencie, gdy kolejna, analizowana decyzja jest już zakwalifikowana do klasy przeciwstawnej. Oczekiwany wynikiem

procesu kwalifikacji jest wzbogacenie wyjściowych przykładów obu przeciwstawnych klas. Wynikową rekomendacją jest zbiór decyzji nieulegający zmianom w kolejnej iteracji.

Na rys. 3 przedstawiono wyniki wyznaczania rekomendacji metodą iteracyjnego uzupełniania wskazanych przez decydenta przykładów decyzji akceptowalnych oraz nieakceptowalnych. Do obliczeń przyjęto takie same wskazania jak w obliczeniach poprzednich. Przyjęto wartość współczynnika regularyzacji $\rho = 0,005$. Ograniczenie liczebności klastrów było uzyskiwane przez ustalenie maksymalnej, dopasowanej do wzorca odległości decyzji od średniej wartości klastra wzorca. Brak zmian rozwiązania zaobserwowano po jedenastu iteracjach.

8. Podsumowanie

Przedstawione sformułowanie problemu oparte jest na zadaniu wyznaczania rekomendacji w procesie decyzyjnego wspomaganie decydenta w wyszukiwaniu informacji w dużych bazach danych (np. w diagnostyce medycznej, automatycznym poszerzaniu zbioru uczącego w zadaniach uczenia sieci neuronowych). Takie ujęcie problemu ułatwia interpretację algorytmu, w tym dobór uniwersalnego słownictwa. Obszar zastosowań można ogólniej określić jako analizę danych z niestandardowego punktu widzenia.

Zasadniczy wynik przedstawionych w artykule propozycji stanowi konstatacja, że współpracę z decydem przy wyznaczaniu ograniczonych rekomendacji można zredukować do jednorazowego wskazania przykładów i kontrprzykładów.

Przedstawiane rozwiązanie bazuje na interpretacji dokonanych przez decydenta wskazań jako przykładowych wzorców dwóch przeciwstawnych klas. Istota proponowanej metody polega na równoległym szeregowaniu decyzji względem wzorców każdej klasy. Rekomendowane klasy uzyskuje się przez wzajemne ograniczanie rekomendowanych sekwencji przez wzorce przeciwstawne. Uzyskane w ten sposób rekomendacje stają się przykładowymi wzorcami w następnej iteracji.

Wymiar generowanej przez przykłady podprzestrzeni cech jest ograniczony przez liczbę wskazywanych przykładów. Wynikający stąd problem polega na możliwości odrzucania decyzji, których wektory cech nie leżą w wygenerowanej przez przykłady podprzestrzeni. Narzędziem pomagającym eliminować takie przypadki jest regularyzacja. Stosowanie zbyt dużych wartości współczynnika regularyzacji prowadzi jednak do rekomendacji środowiskowej.

Przedstawione w artykule wyniki uzyskano przy wykorzystywaniu metody centroidalnej obliczania odległości między klastrami. Większą wrażliwość na

zmiany pojedynczych decyzji można uzyskać, stosując metody bardziej złożone [2].

Ocena preferencji decydenta na podstawie porównywania odległości między klastrami prowadzi do grupowania cech o wartościach skorelowanych. Użyteczność tej metody jest widoczna zwłaszcza w przypadku dysponowania obserwacjami (cechami) niefiltrowanymi z punktu widzenia preferencji decydenta (tzn. kiedy poszczególne cechy nie są wyróżnikami preferencji).

Satysfakcja decydenta-człowieka z uzyskanej rekomendacji jest trudna do przewidzenia. Ulotność problemu jest konsekwencją niewiedzy decydenta. Jest skutkiem ograniczonej zdolności decydenta do ewaluacji dużej liczby decyzji. W zadaniach automatycznego przeszukiwania baz danych (np. w celu poszukiwania podobnych obiektów) można formułować ocenę jakościową proponowanej procedury.

Uzyskanie rozwiązań trywialnych (np. pustych klastrów) może być skutkiem niespójności dokonanych przez decydenta wskazań, wynikających bądź z niedopasowania przestrzeni obserwacji (cech) do sygnalizowanych preferencji decydenta, bądź ze sprzecznych jego wskazań.

Literatura

- [1] BALL G.H., HALL D.J., *Isodata, an Iterative Method of Multivariate Analysis and Pattern Classification*. Proceedings of the IFIPS Congress, 1965.
- [2] KWIATKOWSKI W., *Metody automatycznego rozpoznawania wzorców*. BEL Studio, Warszawa, 2010.
- [3] KWIATKOWSKI W., *Recommendations as a result of decision evaluations based on reference examples*, Teleinformatics Review, No. 1-2, 2019, pp. 3-23.
- [4] KWIATKOWSKI W., *The regularization method in the classification task according to given examples* Teleinformatics Review, No. 3-13, 2019, pp. 3-23.
- [5] KWIATKOWSKI W., *Wykrywanie anomalii bazujące na wskazanych przykładach*. Przegląd Teleinformatyczny, nr 1-2, 2018, s. 3-21.
- [6] MOBASHER B., DAI H., LUO T., NAKAGAWA M., *Improving the Effectiveness of Collaborative Filtering on Anonymous Web Usage Data*. In: Proceedings of the IJCAI 2001, Workshop on Intelligent Techniques for Web Personalization (ITWP01), 2001.

Recommendations as a result of the assessment of preferences on the basis of the indicated examples

ABSTRACT: The problem of determining a decision recommendation according to examples of acceptable decisions and examples of unacceptable decisions indicated by the decision-maker is considered in the paper. The decision-maker's examples are the foundation for assessing his preferences. The essence of the presented solution consists in determining the preferences of the decision-maker as a cluster designated by supplementing the indicated examples. The paper proposes a procedure of successive approximations based on the classification task according to given examples.

KEYWORDS: recommendation, preferences, data mining, classification, clustering

Praca wpłynęła do redakcji: 18.10.2021 r.

Analysis and classification of chosen social engineering methods in cybersecurity¹

Monika Olchowik

Institute of Teleinformatics and Cybersecurity, Faculty of Cybernetics, MUT,
ul. gen. Sylwestra Kaliskiego 2, 00-908 Warsaw, Poland
monika.olchowik.pl@gmail.com

ABSTRACT: Cyberthreat landscape is everchanging and dynamically evolving. Tools, techniques and software are getting more and more intricate. In contrast social engineering methods have been used in various attacks long before computers have been created, yet they are as useful as before, even in cyberspace. Social engineering attacks are quite often successfully used by conmen and hackers, and as such are a constant part of cyberthreat landscape. In order to detect and prevent the usage of aforementioned techniques greater understanding and systematisation of the process is need. In this paper a classification of chosen social engineering methods has been proposed. The classification is based on Kevin Mitnick's Social Engineering Cycle. This classification allows for creation of attack patterns and could be used as a basis for a social engineering attack matrix. Moreover, the paper presents a collection of different methods used in each of the stages of the cycle, describes them and provides examples of their usage.

KEYWORDS: Social Engineering, Social Engineering Cycle, Social Engineering Methods, Social Engineering Methods Classification

1. Introduction

Digital revolution is happening even faster with each passing year. It influences not only how business is conducted, but it also changes all aspects of life. From communication to transportation, technology is ever-present and everchanging. As such the cybersecurity threat landscape is constantly evolving. With each passing year new vulnerabilities, methods and vectors of attacks are discovered and used. In response organizations develop new procedures, patch

¹ I dedicate this article to the memory of my promoter, Dr. Eng. Zbigniew Suski

and upgrade their software. Finally, they constantly invest in even more sophisticated systems for better detection and defence.

Yet no matter the technology used, there is one unchanging component – the technology is made for and used by people. This allows to exploit the oldest vulnerability, that is the human nature. Such attacks might be performed by social engineering. This process can be defined as “manipulating people, by deception, into giving out information, or performing an action [1]”.

Even up to 34% of organizations consider careless/unaware employees as the biggest vulnerability [2]. Employee weakness is considered to be responsible for 20% of confirmed breaches by 2020 [3]. In order to minimise this risk, it is important to understand social engineering and its methods.

2. The Social Engineering Models

One of the most well-known social engineering models is Mitnick’s Social Engineering Cycle. It consists of four stages: Research, Developing Rapport and Trust, Exploiting Trust, and Utilizing Information [4]. The cycle is shown on Figure 1.



Figure 1. The Social Engineering Cycle

During the initial part called Research attacker tries to gather as much information as possible about a target. In this phase it is important to acquire data about the people, the company, but also about the targeted social group.

In the second phase the attacker uses the information gathered to develop the trust between himself and the victim. There are many approaches to this phase, all of them deeply rooted in psychology. [4] The attacker could assume position of authority over the victim or try to be friendly and approachable. It is also possible for this phase to take longer. The attacker can slowly introduce

themselves into the surroundings of the victim. If performed successfully the victim will recognize the attacker as the part of the chosen community, or the company. This phase creates the base of the success of a social engineering attack. The victim is more likely to be vulnerable to an attack if they believe the attacker.

The third step is the key of the social engineering attack. Its goal is to convince the victim to perform an action or to give out information. The attacker could simply ask the victim about the needed data or ask for help. However, in a process called the reverse sting the victim might be manipulated to request help from the attacker and by following the instructions of the attacker realising his current goal [4].

The fourth phase called Utilise Information is dependent on the attacker acquiring all the necessary information. In this case the goal of the attacker is achieved, and no more actions are necessary. Otherwise the attacker continues and performs the phase one again.

It is worth noting that there are other social engineering models. An example of this is the Social Engineering Attack Framework. It has been inspired by the Social Engineering Cycle [5].

The Social Engineering Attack Framework is divided into six phases: Attack Formation, Information Gathering, Preparation, Relationship Development, Relationship Exploit, and Debrief. Each of the phases consists of individual steps [5]. Individual steps of this model are shown on Figure 2. It is of note that “Goal Satisfaction” does not belong to any of the phases.

Individual steps of The Social Engineering Attack Framework (Figure 2) are coloured according to their phase: brown (Attack Formation), orange (Information Gathering), violet (Preparation), dark blue (Relationship Development), light blue (Relationship Exploit), green (Debrief), teal (not belonging to any phase).

It is possible to notice correlation between phases of those two models. The Research phase has been split into Attack Formation, Information Gathering and Preparation. Developing Rapport and Trust is a direct equivalent of Relationship Development. Exploiting Trust is covered by Relationship Exploit. Utilizing Information corresponds to two actions: Transition and Goal Satisfaction.

The difference in the two aforementioned models is the maintenance, which is not explicitly shown in the Social Engineering Cycle. However, it is possible that such actions are a part of the Exploiting Trust phase. In this action the attacker is trying to calm the victim and appease their emotional state. The goal of this step is to ensure that the victim does not feel as if they have been attacked [5] and thus do not perform any actions such as changing passwords or reporting the incident.

Due to the similarities between the models in this article methods will be classified by their role in accordance with the Social Engineering Cycle, as it is the most often used model to describe social engineering attacks.

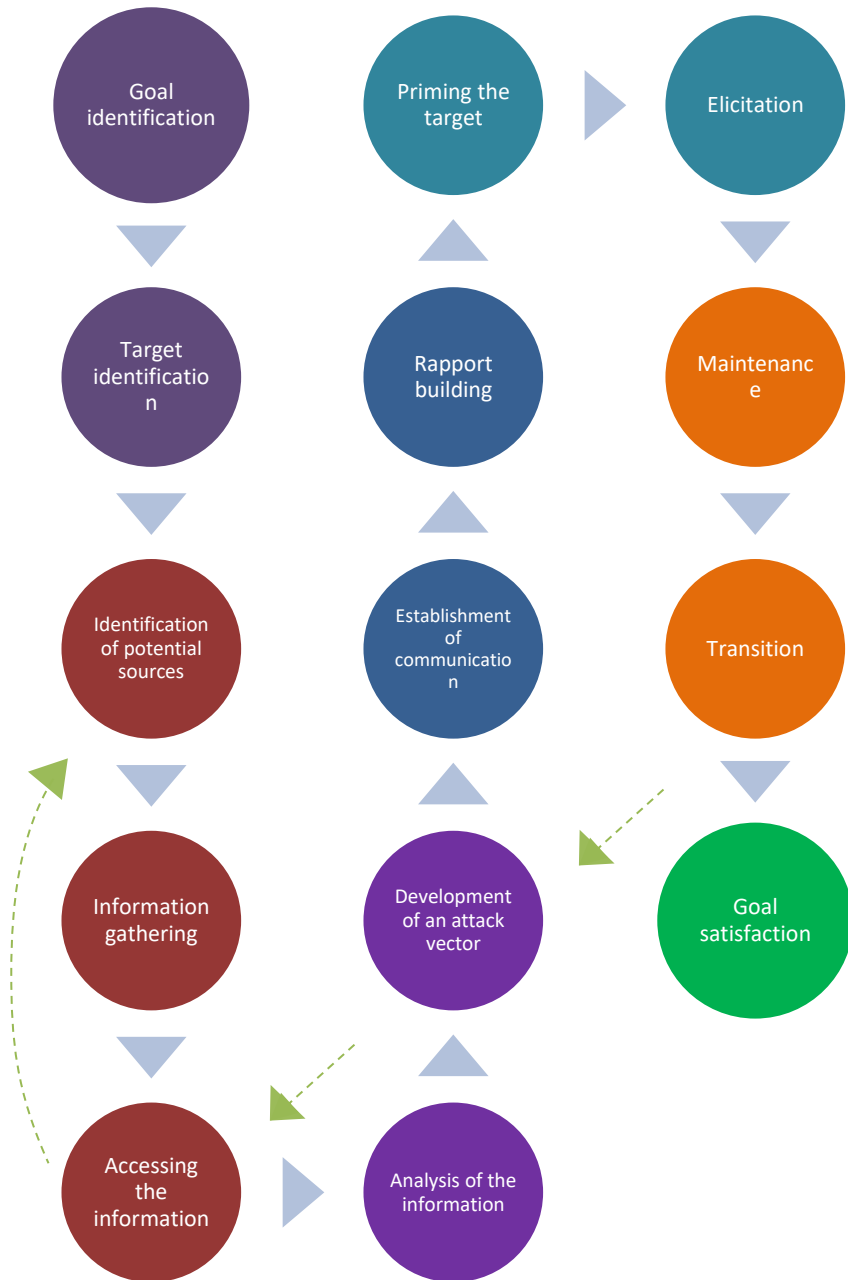


Figure 2. The Social Engineering Framework

3. Classification of methods based on the Social Engineering Cycle

3.1. Basis of the classification

Methods analysed in this paper have been chosen based on the prevalence of their usage and their applicability to various scenarios [1, 4, 8]. Thus, it is possible to attempt to classify these general-purpose methods based on the stage of the social engineering attack, in which they would be used. Each of the phases has a common goal and thus the basis of the proposed classification is the result achieved by the usage of a given method. However, the last phase of the cycle – Utilising Information is a notable exception as its goals are highly case specific and thus this phase has not been included as a part of the categorization.

The proposed classification of different sociotechnical methods introduces following categories:

- Research,
 - Passive information gathering,
 - Active information gathering,
- Developing Rapport,
- Trust and Exploiting trust.

Moreover, the methods used in the Research phase have been split in two categories in order to differentiate between the Passive and Active Information Gathering.

3.2. Research

During the research phase the goal of an attacker is to find as much information about his target as possible. The methods used in this phase can be categorized in analogous way as the scanning methods. The distinction between the methods is whether the target can learn or suspect that the information is being gathered.

3.2.1. Passive information gathering

Passive information gathering is centred on the idea of acquiring data without the victim being aware about it. As such open-source intelligence (OSINT) tools are used at this stage. The OSINT Framework [6] offers a collection of websites and tools used in such processes. Among the important

sources used during a reconnaissance are:

- search engines,
- company website,
- unsecured cloud storage,
- people databases,
- social networks and dating sites,
- tools for analysis of social networks,
- databases of leaks from online communicators,
- online registries and records,
- forums and blogs,
- files and their metadata.

Social media is a source especially rich in personal data. It is an ideal source of data for social engineers who would like to commit identity theft. Sources like Facebook, Instagram, Twitter or LinkedIn allow the attacker to learn about the mannerisms of a victim, important events from their life and other personal information [7]. All of this is willingly shared by the target. For example, knowing that an important executive in a company is on vacation on the other side of the world, a social engineer might safely assume the identity of a person representing the executive, e.g. a new assistant.

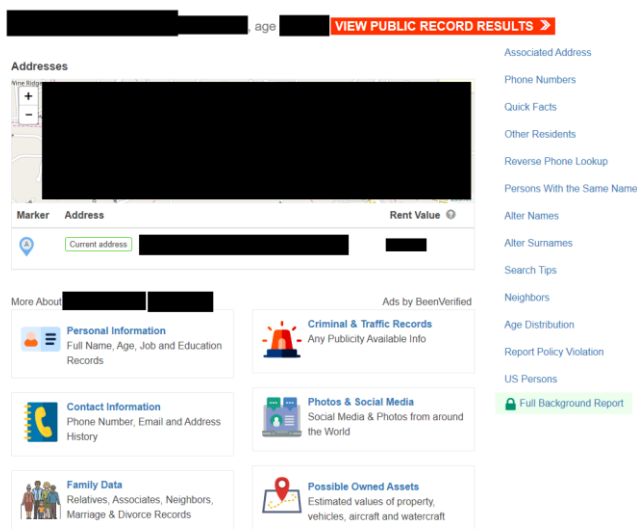


Figure 3. An example people database [9]

It is worth noting that people databases contain collections of personal information about many people. The knowledge of the person's name and surname is required and usually sufficient to use them. Additionally, some of them offer paid services. An example is shown on Figure 3. The attacker can also use highly specific tools such as Maltego. This tool helps with the information gathering, automating some of the actions required, and as such makes the process easier for the attacker [8].

3.2.2. Active information gathering

Active information gathering requires a degree of sophistication. An attacker risks that a victim will realize that they are being targeted and in turn will be on high alert. It is also possible that the attacker is exposing himself to a greater danger. Most notable example of active information gathering is people watching. If done incorrectly the attacker can be spotted and engaged by the victim. If done correctly however, one might be able to learn the occupation, personal details, habits and interests of the victim. A cybersecurity researcher by the name of Johnny Long was able to spot and determine the identity of a government agent at the airport [10]. During people watching it is possible to find a pattern in corporate attire and analyse badges. A good replica of a badge that supposedly stopped working might be enough for a security team to let a person inside of a building [10].

If a person watched by a social engineer uses a computer or a phone, the attacker might try shoulder surfing. This action is a simple act of watching the victim's device over their shoulder in order to gather important information that is being shown on the screen [10] or typed [11]. In some cases, it might be used to obtain the unsuspecting victim's password [4]. If two co-workers are chatting, eavesdropping might be an additional source of relevant information. The social engineer can learn through this e.g. their nicknames [12], or the projects the employee is participating in.

A similarly useful, simple and widely applicable technique is dumpster diving. An attacker might need to enter a property to engage in this activity. However, it can yield surprisingly useful results. It is possible to find personal information, corporate classified data, and financial information this way [10]. Vehicle watching is a different, supplementary action. It is based on looking at documents and receipts left unattended by a victim in their car.

A social engineer requires good observational skills. He not only needs to analyse badges and various documents, but also has to understand the inner workings of the company's building. The attacker should be aware of the layout of the building, its entries and exits, guards and other security staff, requirements for entry, and location of keypads. This information might be needed as shown in

the example of “Big Blue Pest Control” [8]. In this case social engineers were rejected by the security guards in the lobby. However, they remembered their names. Then they spotted the entry used by smokers, with lower security than the main entry. By pretending to be inspecting the building they successfully followed others into the building and then inside of it into an elevator. Later they claimed that they were referred by the security guard whose name they remembered. This was enough for them to enter the mailroom.

A different approach to active information gathering is the usage of all sorts of personality quizzes [11]. In this case the attacker disguises his true goal. Quite often a victim is enticed by the potential reward – a funny and memorable result. If the quiz is memorable enough the victim can share it with their friends. One of the most straightforward examples spotted in the wild was a questionnaire called “What does your password say about you?” [11]. It was created as a Facebook app, and contained questions about length of the password, and its complexity. However other quizzes e.g. testing compatibility between people or determining the job that would suit the user, can be used to determine the personality profile of the victim or provide enough information to impersonate a person. A live survey about Easter candy eating habits conducted by InfoSecurity Europe in 2006 shown that in 81% the researchers gathered enough of information to attempt an identity theft [12].

3.3. Developing Rapport and Trust

In this phase the Social Engineer engages the victim. If the communication is initiated in person, non-verbal communication is of high importance. As such, the verbal statements have to be corresponding with the message sent by micro and macroexpressions, voice, body language and gestures [13]. These elements are necessary for a successful elicitation to occur, that is to learn useful details about the target in the course of a seemingly normal and friendly conversation [13].

However, what is more surprising is that non-verbal communication might be applied even in a phishing attack. Emoticons might be used to convey openness and friendliness when trying to befriend the victim on social media platforms and dating sites. A different technique of non-verbal communication is called framing [13], that is using pictures, fonts, formatting etc. so that the victim thinks that the malicious content is part of a legitimate site or communication.

In general, common behaviours of social engineers in this phase include:

- assumption of authority,
- showing confidence,
- impersonation and pretexting:

- name dropping,
- insider knowledge including usage of the lingo,
- usage of leaked data,
- excuses;
- befriending the target:
 - flattery and validation,
 - flirting,
 - asking open ended questions,
 - quid pro quo;
- using sympathetic / guilt-inducing themes;
- elicitation;
- setting time constraints or creating artificial scarcity of items / information.

All the aforementioned behaviours might be spotted in both personal interactions, as well as communication over the Internet. The example of the latter are extortion campaigns. In the first five months of 2019, 300 million of such emails have been stopped by Symantec [14]. One of the most popular type is sextortion. The attacker informs the victim that he knows what their current password is. To establish trust attacker uses leaked data hoping that the target hasn't changed their password since then. Then the social engineer impersonates a hacker claiming that he hacked a website and infected the victim's computer, essentially gaining access to the victim's camera. Through that, he collected videos of illicit nature. As such the attacker assumes the position of authority and threatens the target that the video will be made public if they won't comply with the request. Campaigns following this simple scheme have been popular in 2019 [14][15] and 2020 [16]. Even though these attacks focus on intimidating the target, trust is still required. Otherwise, the victim would not believe in the possible negative consequences.

Sympathy might be dangerous as well. Assuming an identity of a desperate jobseeker that unfortunately spilled coffee over his resume is enough to elicit sympathetic feelings. Even more, it might be just enough to convince a receptionist to plug a USB stick with malicious files in and open them [17].

Curiosity and the need to reciprocate a favour might be also used against the target. Quid pro quo is a method in which the attacker offers something of perceived similar value. It can be as simple as providing feedback in a conversation e.g. information about the created persona [13], or as complicated as offering free services in exchange for access to information [18].

A different way of establishing trust is elicitation, especially if combined with impersonation. A hacker group UG-NAZI has extracted credit card data thanks to tech support performing a password reset on database admin account.

The tech support tried to verify if the attacker on the phone is really the database admin, by asking different questions that the attacker knew the answers to. As such by sharing information, he gained trust of the victim [13].

3.4. Exploiting trust

Social engineer in this phase tries to convince the target to provide information or perform an action that will satisfy the current goal of the attacker. In this phase following methods are used:

- questions
- baits
- threats and negative consequences for not fulfilling requests
- reciprocation
- reverse sting
- distraction

Exploiting trust might be as simple as nicely asking the victim to perform an action, as it was shown in the example of a desperate jobseeker and a helpful receptionist. However, it isn't always enough. In extortions attacker bases his message on threats e.g. of publishing videos, or even infecting the victim's family with SARS-COV-2 [18].

On the other hand, a phishing campaign "2011 Recruitment Plan" targeting RSA shows the importance of baiting. The attack targeted a small group of employees. However, it was crafted well enough to look like a legitimate message. The email had a malicious Excel spreadsheet attached. For one of the targets the information about a possible recruitment plan was interesting enough for them to open it [7].

Similar effect has the usage of authority and giving orders or by reciprocation. Helping a person makes them more likely to comply with requests. Especially impressive effects can be created if it is used with a reverse sting approach. That is an attack in which the attacker manipulates the victim into turning to the attacker for help. One of the examples of this is a final stage of an attack on a new employee in a company. The attacker is aware of the name, surname and phone number of the target. Social engineer introduces himself as a member of information security team and offers to help the victim with all the intricacies of cybersecurity policies. After walking through some of them he asks about the complexity of the current password. As the new employee did not have a complex password the attacker proposed a change to a new password that they have created together [4]. The victim was grateful for help and provided each answer to the fullest capacity. However, they have unknowingly endangered the

company. Another example of this is a story of an unnamed social engineer mentioned by Kevin Mitnick. First, the attacker called a publicly known phone number for a sheriff's station. He introduced himself as a police officer and claimed that he tried calling a different number and implied that he must have made a mistake. In turn, the local police officer provided him with the correct phone number for internal use. Thus, the attacker has received the information he wanted even without asking for it [4].

A different approach is used if the goal is to make sure that a person does not perform an action. Then a distraction might be just enough. The Whurley's exploit shows that a story about a work colleague not returning the money and thus the attacker lacking money to take a date out might be distracting enough not only to bypass the security check, but also being given dating advice and money for lunch [19].

4. Summary

The social engineering cycle is an important tool that can be used to understand and describe social engineering attacks. Even a seemingly simple attack such as phishing might be described using a single cycle. Usually attacker has to perform research, impersonate a person or an organization in order to establish trust and legitimacy, and in the end ask, threaten or convey in any other way that the victim should open a malware-ridden attachment or a malicious link.

The simplicity of this model makes remembering and understanding it easy. Thus, it allows for an assignment of social engineering methods and techniques to the phases of the cycle. The end result of such classification could be used for creation of a social engineering attack matrix. A formalized attack matrix could allow for faster detection or easier analysis of the attack.

References

- [1] MANN I., *Hacking the Human: Social Engineering Techniques and Security Countermeasures*. Aldershot: Gower, 2008.
- [2] *Is cybersecurity about more than protection?* EY Global Information Security Survey 2018-19. EYGM Limited, 2018.
- [3] *How does security evolve from bolted on to built-in? Bridging the relationship gap to build a business aligned security program*. EY Global Information Security Survey 2020. EYGM Limited, 2020.
- [4] MITNICK K.D., SIMON W.L., *The Art of Deception: Controlling the Human Elements of Security*. Wiley Publishing, Indianapolis, 2002.

- [5] MOUTON F., MALAN M.M., LEENEN L., VENTER H.S., *Social Engineering Attack Framework*. In: Information Security for South Africa, Conference Paper, South Africa, Johannesburg, 2014.
- [6] *OSINT Framework*, 5.2.2020. [Online]. Available: <https://osintframework.com/>. [Accessed 16.4.2020].
- [7] ALEXANDER M., *Methods for Understanding and Reducing Social Engineering Attacks*. The Sans Institute, 2016.
- [8] HADNAGY C., *Social Engineering. The Science of Human Hacking*. Wiley, Indianapolis, 2018.
- [9] *Public Records Encyclopedia*, ClusterMaps.com, [Online]. Available: <https://clustrmaps.com/>. [Accessed 16.4.2020].
- [10] LONG J., *No Tech Hacking. A Guide to Social Engineering., Dumpster Diving and Shoulder Surfing*. Elsevier, Burlington, 2008.
- [11] BROWER J., *Which Disney© Princess are YOU? (Web 2.0) Social Engineering on Social Networks*. The SANS Institute, 2010.
- [12] MANJAK M., *Social Engineering Your Employees to Information Security*. SANS Institute, 2006.
- [13] HADNAGY C., EKMAN P., *Unmasking the Social Engineer: The human element of security*. Wiley, Indianapolis, 2014.
- [14] *Symantec Enterprise Blog*, Symantec, 30 July 2019. [Online]. Available: <https://symantec-blogs.broadcom.com/blogs/threat-intelligence/email-extortion-scams>. [Accessed 16.04.2020].
- [15] *Naked Security*, Sophos. 17.12.2019. [Online]. Available: <https://nakedsecurity.sophos.com/2019/12/17/dont-fall-for-this-porn-scam-even-if-your-passwords-in-the-subject/>. [Accessed 17.4.2020].
- [16] ABRAMS L., *Bleeping Computer*, *Bleeping Computer*, 9.4.2020. [Online]. Available: <https://www.bleepingcomputer.com/news/security/large-email-extortion-campaign-underway-dont-panic/>. [Accessed 17.4.2020].
- [17] HADNAGY C., *Social Engineering: The Art of Human Hacking*, Wiley, Indianapolis, 2011.
- [18] *Naked Security*, Sophos. 19.3.2020. [Online]. Available: <https://nakedsecurity.sophos.com/2020/03/19/dirty-little-secret-extortion-email-threatens-to-give-your-family-coronavirus/>. [Accessed 17.4.2020].
- [19] SIMON W.L., MITNICK K.D., *The Art of Intrusion: The Real Stories Behind the Exploits of Hackers, Intruders and Deceivers*. Wiley, Indianapolis, 2005.

Analiza i klasyfikacja wybranych metod socjotechnicznych stosowanych w cyberbezpieczeństwie

STRESZCZENIE: Cyberbezpieczeństwo jest dziedziną dynamicznie się zmieniającą. Narzędzia, techniki, oprogramowanie są ciągle rozwijane i stają się coraz bardziej złożone. W przeciwieństwie do nich, metody socjotechniczne używane były od wielu lat i nadal są aktualne, nawet gdy wykorzystywane są w cyberprzestrzeni. Ataki socjotechniczne są często przeprowadzane z sukcesem przez oszustów oraz hackerów. Niezmiennie pozostają one zagrożeniem. W celu wykrycia i przeciwdziałania takim technikom konieczne jest zrozumienie i usystematyzowanie procesu ich wykorzystania. W niniejszym artykule zaproponowano klasyfikację wybranych metod socjotechnicznych opartą o Cykl Socjologiczny Kevina Mitnicka. Klasyfikacja pozwala na tworzenie wzorców ataku oraz może zostać użyta w celu stworzenia matrycy ataków socjotechnicznych. W artykule przedstawiono również zbiór różnych metod socjotechnicznych używanych w każdym z etapów cyklu, opisano je oraz podano przykłady ich zastosowania.

SŁOWA KLUCZOWE: socjotechnika, cykl socjotechniczny, metody socjotechniczne, klasyfikacja metod socjotechnicznych

Received by the editorial staff on: 1.06.2021

Artykuł poświęcam pamięci mojego promotora dr inż. Zbigniewa Suskiego.

Ocena dokładności i powtarzalności pomiaru parametru RSSI w systemach BLE

Artur ARCIUCH

Instytut Teleinformatyki i Cyberbezpieczeństwa, Wydział Cybernetyki WAT,
ul. gen. Sylwestra Kaliskiego 2, 00-908 Warszawa
artur.arciuch@wat.edu.pl

STRESZCZENIE: Systemy kontroli dostępu są powszechnie stosowane w organizacjach. Do kontroli dostępu wykorzystywane są specjalizowane karty oraz czytniki. Wykorzystanie w takich systemach kart oraz czytników wyposażonych w interfejs BLE eliminuje konieczność zbliżenia karty do czytnika w celu uzyskania dostępu do zasobu organizacji. W artykule oceniono użyteczność urządzeń z interfejsem BLE w systemach automatycznej kontroli dostępu, opierającej się na dokładności i powtarzalności pomiaru parametru RSSI.

SŁOWA KLUCZOWE: BLE, RSSI, pomiar odległości

1. Wprowadzenie

Urządzenia z interfejsami do komunikacji bezprzewodowej są obecnie powszechnie wykorzystywane. Przykładem popularnego interfejsu do komunikacji bezprzewodowej na bliskie odległości: 0-10 m¹ jest Bluetooth Low Energy (BLE), w który są wyposażone urządzenia IoT (ang. Internet of Things), w tym urządzenia typu mikrokontroler, czy tzw. znaczniki Bluetooth (ang. beacon). Charakteryzują się one: małą masą i rozmiarami oraz stosunkowo długim czasem pracy i niskim poborem mocy. Są one stosowane do lokalizowania obiektów wewnątrz budynków [1], [2], [9], [10]. Urządzenia typu sensor (nasłuchujące), rozmieszczone w znanej lokalizacji, dokonują pomiaru wartości parametru RSSI (Received Signal Strength Indicator / Index) sygnału generowanego przez urządzenia typu marker (znacznik). Wartość RSSI jest

¹ Klasa mocy 2 (2,5 mW).

przeliczana na odległość markera od sensora. Lokalizacji obiektu dokonuje się, wykorzystując np. metodę triangulacji [1], [2], [4], przy zastosowaniu m.in. urządzeń z interfejsem Bluetooth (BLE w szczególności): Bluetooth beacon, iBeacon, mikrokontroler, smartfon itp. [1], [2], [4], [7], [8].

Systemy nadzoru dostępu do pomieszczeń, te obecnie najczęściej spotykane, wymagają instalacji specjalistycznych urządzeń typu karty i czytniki kart. Użytkownik w celu wejścia do pomieszczenia musi zbliżyć kartę do czytnika na bardzo bliską odległość. Ta czynność wymaga na ogół użycia jednej dłoni. Ponadto system nadzoru może zarejestrować jedynie informację o dostępie do czytnika.

Prezentowana w artykule propozycja polega na wykorzystaniu w systemie nadzoru dostępu kart nie wymagających bliskiego kontaktu z czytnikami. Użytkownicy wyposażeni są w markery (beacony), pełniące rolę kart, poruszają się w obrębie budynku, w którym rozmieszczone są sensory wykrywające markery i wykonujące pomiary odległości. Sensor przesyła do serwera informacje zawierające m.in. identyfikatory wykrytych markerów oraz odpowiadające im wartości parametrów RSSI. Oprogramowanie serwera, na podstawie informacji uzyskanych z sensorów, dokonuje śledzenia oraz analizy zmian położenia markerów w celu automatycznego przydzielania albo zabrania dostępu do pomieszczeń albo stref, poprzez np. sterowanie zamkami drzwi.

Dokonana w artykule ocena użyteczności urządzeń z interfejsem BLE w systemach automatycznej kontroli dostępu opiera się na ocenie dokładności i powtarzalności pomiaru parametru RSSI.

2. Metoda wyznaczania odległości

Przedmiotem prezentowanego w artykule badania był problem wyznaczania odległości pomiędzy urządzeniami z interfejsem BLE w pomieszczeniach budynku (o powierzchni mniejszej niż 25 m²). W szczególności celem eksperymentu było zbadanie dokładności i powtarzalności pomiaru wartości parametru RSSI w zależności od użytych egzemplarzy markerów tego samego typu i egzemplarzy sensorów tego samego typu. W przeprowadzonym eksperymencie dokonywano pomiarów odległości między sensorami beetle (DFRobot Beetle ESP32 v2.0²) a markerami iNode Beacon³.

Metoda wyznaczenia odległości między sensorem a markerem polegała na pomiarze przez sensor wartości parametru RSSI sygnału markera, a następnie

² <https://www.dfrobot.com/product-1798.html> (30.09.2021)

³ <https://inode.pl/iNode-Beacon,p,17> (30.09.2021)

wyznaczeniu odległości w funkcji RSSI. Parametr RSSI jest określany jako wskaźnik poziomu sygnału odbieranego przez urządzenie bezprzewodowe. Na potrzeby obliczania zależności RSSI od odległości między antenami markera i sensora przyjęto model wykładniczy [2], [5] opisany równaniem (1):

$$RSSI(d) = RSSI(d_0) - n \cdot 10 \cdot \log_{10} \left(\frac{d}{d_0} \right), \quad (1)$$

gdzie:

d_0 – odległość odniesienia,

$RSSI(d_0)$ – wartość RSSI zmierzona w odległości d_0 ,

d – odległość pomiaru,

$RSSI(d)$ – wartość RSSI zmierzona w odległości d ,

n – odległościowy współczynnik tłumienia (PLE – ang. path loss exponent).

Parametr n zależy od propagacji sygnału w danym otoczeniu. Przyjmuje się [6], że w budynkach mieszkalnych i biurowych $n \in \langle 1,6,6 \rangle$.

Po przekształceniu zależności (1) uzyskuje się następującą zależność do wyznaczenia wartości odległości [2]:

$$d = d_0 \cdot 10^{\frac{RSSI(d_0) - RSSI(d)}{10 \cdot n}}. \quad (2)$$

Zależność (1) można zapisać w następującej postaci: $y = A + Bx$, gdzie A charakteryzuje wartość parametru RSSI dla odległości d_0 , a B – wartość n określającą warunki propagacji sygnału w środowisku, w jakim dokonywano pomiarów (tzw. PLE). Po wykorzystaniu tych oznaczeń zależności (1) oraz (2) przyjmą odpowiednio postać:

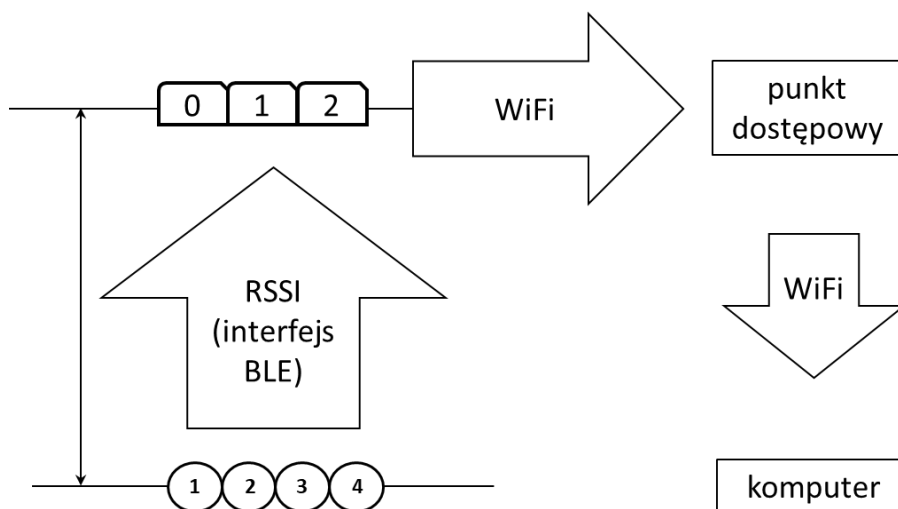
$$RSSI(d) = A + B \cdot 10 \cdot \log_{10} \left(\frac{d}{d_0} \right), \quad (3)$$

$$d = d_0 \cdot 10^{\frac{RSSI(d) - A}{10 \cdot B}}. \quad (4)$$

Schemat stanowiska pomiarowego zaprezentowano na rysunku 1. Stanowisko składało się z trzech sensorów beetle (0, 1, 2 – oznaczone przez trzy prostokąty z liczbami w środku), czterech markerów iNode (1, 2, 3, 4 – oznaczone przez cztery okręgi z liczbami w środku), punktu dostępowego oraz komputera. Utworzono grupę złożoną z trzech sensorów oraz grupę złożoną z czterech

markerów. Elementy grup ułożono wzdłuż dwóch linii, obok siebie, a grupy względem siebie ułożono równolegle. Odległość była ustalana między grupami. Z uwagi na niewielkie rozmiary elementów uznano, że takie ułożenie elementów nie będzie miało istotnego wpływu na dokładność pomiaru.

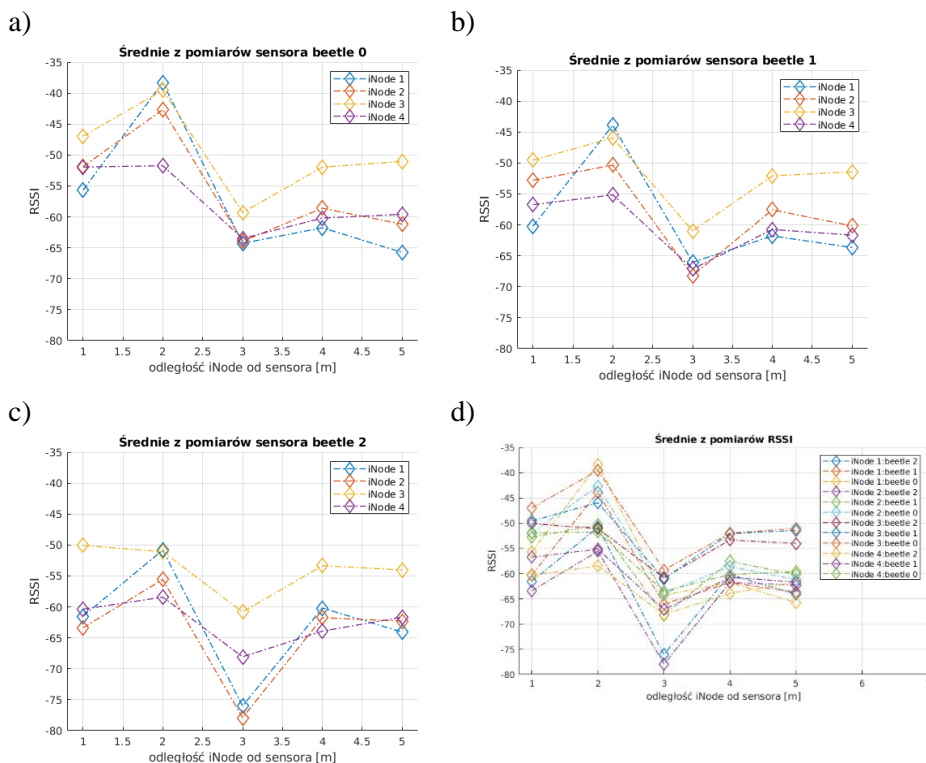
Pomiar polegał na ułożeniu grup w określonej odległości względem siebie. Następnie sensory dokonywały pomiaru wartości parametru RSSI generowanego przez markery. Częstotliwość pomiaru dla każdego sensora wynosiła 0,1 Hz. Wynik pomiaru uzyskany z interfejsu BLE sensora, opatrzony identyfikatorem sensora i identyfikatorem markera, był przesyłany za pośrednictwem sieci WiFi do aplikacji działającej na komputerze, gdzie był rejestrowany. Dla zadanej odległości wykonywano dwieście pomiarów wartości parametru RSSI przez każdy sensor dla każdego markera. Pomiarzy przeprowadzono dla odległości 1, 2, 3, 4 oraz 5 m, między grupami.



Rys. 1. Stanowisko pomiarowe

3. Wyniki pomiarów

Wyniki pomiarów zaprezentowano na wykresach (rysunek 2). Wykresy 2a)-2c) pokazują uśrednione wyniki pomiarów wartości parametru RSSI uzyskiwane przez dany sensora, a wykres 2d) – wyniki pomiarów dla wszystkich sensorów.



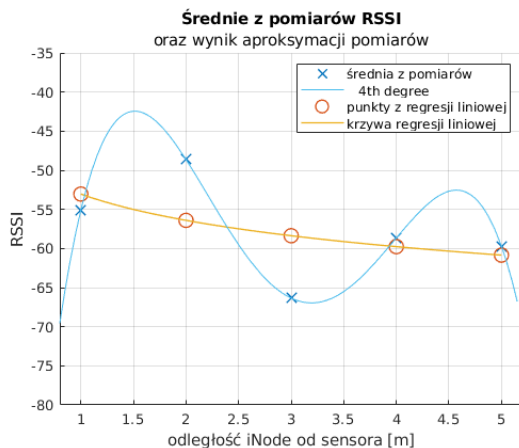
Rys. 2. Uśrednione wyniki pomiarów wartości parametru RSSI uzyskiwane przez sensory: beetle 0 (a), beetle 1 (b), beetle 2 (c), wszystkie (d)

Wyniki pomiarów uśredniono oraz poddano aproksymacji metodą regresji liniowej dla zależności (3) oraz przy wykorzystaniu narzędzia Matlab: basic fitting dla aproksymacji wielomianem czwartego stopnia. Uzyskane zależności opisano, odpowiednio, funkcjami (4) i (5):

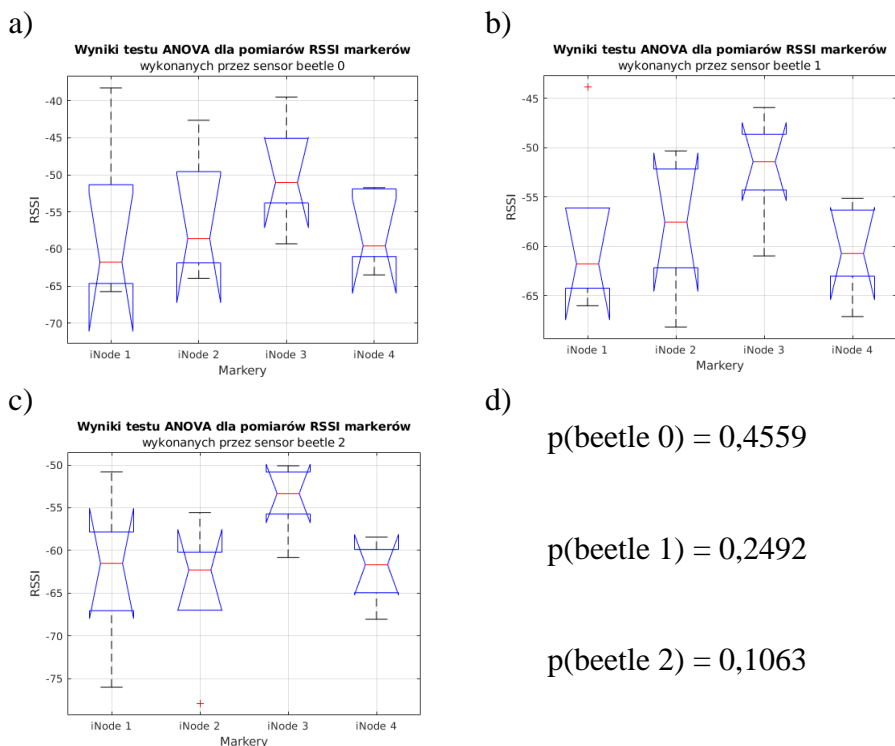
$$RSSI(d) = -53.0397 - 1.1171 \cdot 10 \cdot \log_{10} \cdot \left(\frac{d}{d_0}\right), \quad (4)$$

$$RSSI(d) = -3.491d^4 + 43.18d^3 - 183d^2 + 308.3d - 219.2. \quad (5)$$

Wyniki uśrednionych pomiarów oraz przebiegi funkcji – wyniki aproksymacji pomiarów pokazano na rysunku 3. Na rysunku 4 zaprezentowano rezultaty porównania zgodności wyników uzyskiwanych przez różne sensory dla różnych markerów. Do porównania zastosowano wariant metody ANOVA (ang. Analysis of Variance) zwany jednoczynnikową analizą wariancji. Metoda służy do badania obserwacji, które zależą od jednego czynnika (w rozważanym przypadku tym czynnikiem jest wartość mierzonego parametru RSSI). Wynikiem działania metody jest prawdopodobieństwo, z jakim mierzona wartość RSSI jest powodem różnic między obserwowanymi średnimi grupowymi.



Rys. 3. Wyniki aproksymacji uśrednionych wyników pomiarów wartości RSSI



Rys. 4. Wyniki jednoczynnikowego testu ANOVA dla wartości parametru RSSI mierzonych przez sensor: beetle 0 (a), beetle 1 (b), beetle 2 (c) oraz wyniki działania metody ANOVA (d)

Uzyskane rezultaty – wartości prawdopodobieństw dla sensorów beetle0, beetle1 oraz beetle2, pokazane na rysunku 4d, wskazują na różnice w wynikach pomiarów dla poszczególnych sensorów.

Tab. 1. Porównanie wyników pomiarów parametru RSSI dla par markerów – wyniki były rejestrowane przez sensor beetle 0

<i>A</i>	<i>B</i>	<i>p</i>
iNode 1	iNode 2	0,9917
iNode 1	iNode 3	0,5170
iNode 1	iNode 4	1,0000
iNode 2	iNode 3	0,6850
iNode 2	iNode 4	0,9872
iNode 3	iNode 4	0,4908

Tab. 2. Porównanie wyników pomiarów parametru RSSI dla par markerów – wyniki były rejestrowane przez sensor beetle 1

<i>A</i>	<i>B</i>	<i>p</i>
iNode 1	iNode 2	0,9895
iNode 1	iNode 3	0,3636
iNode 1	iNode 4	0,9924
iNode 2	iNode 3	0,5311
iNode 2	iNode 4	0,9356
iNode 3	iNode 4	0,2436

Tab. 3. Porównanie wyników pomiarów parametru RSSI dla par markerów – wyniki były rejestrowane przez sensor beetle 2

<i>A</i>	<i>B</i>	<i>p</i>
iNode 1	iNode 2	0,9794
iNode 1	iNode 3	0,2180
iNode 1	iNode 4	1,0000
iNode 2	iNode 3	0,1133
iNode 2	iNode 4	0,9783
iNode 3	iNode 4	0,2205

W tabelach 1-3 pokazano rezultaty porównania wyników pomiarów parametru RSSI dla par markerów (kolumny *A* i *B* opisywały parę markerów, kolumna *p* – wartość określająca podobieństwo wyników), uzyskane za pomocą metody

jednoczynnikowej analizy wariancji – wyniki były rejestrowane odpowiednio przez sensory: beetle 0, beetle 1 i beetle 2. Wartości w kolumnie p bliskie 1 oznaczają zbliżone wartości pomiarów parametru RSSI.

4. Wnioski

Stanowisko pomiarowe może stanowić szkielet systemu automatycznej kontroli dostępu. Z wyników pomiarów wynika, że parametr RSSI dla wykorzystanych do badań sensorów i markerów przyjmuje wartości z przedziału (-80, -40), a wartości średnie z pomiarów – (-62, -53). Częstotliwość pomiarów dokonywanych przez sensor wynosiła 0,1 Hz. Wiersze szarego koloru w tabelach 1-3 oznaczają, że pomiary wartości parametru RSSI dla zaznaczonych par markerów są bardzo zbliżone. Takie zbliżone wyniki zostały uzyskane przez wszystkie sensory dla wszystkich markerów.

Wyniki pomiarów wskazują, że aby zwiększyć dokładność wyznaczenia odległości na podstawie pomiaru wartości parametrów RSSI, z wykorzystaniem użytych w eksperymentach urządzeń, należy dokonać kalibracji każdej pary sensor – marker niezależnie, w fazie instalacji systemu.

Literatura

- [1] CHAI S., AN R., DU Z., *An Indoor Positioning Algorithm Using Bluetooth Low Energy RSSI*. International Conference on Advanced Material Science and Environmental Engineering (AMSEE 2016), 2016.
- [2] THALJAOU A., VAL T., NASRI N., BRUBLIN D., *BLE localization using RSSI measurements and iRingLA*. 2015 IEEE International Conference on Industrial Technology (ICIT), Seville, Spain, 2015, pp. 2178-2183.
- [3] JIN R., CHE Z., XU H., WANG Z., WANG L., *An rssi-based localization algorithm for outliers suppression in wireless sensor networks*. Wireless Networks, vol. 21, no. 8, 2015, pp. 2561-2569.
- [4] WANG Y., YANG X., ZHAO Y., LIU Y., CUTHBERT L., *Bluetooth positioning using RSSI and triangulation methods*. 2013 IEEE 10th Consumer Communications and Networking Conference (CCNC), Las Vegas, NV, USA, 2013, pp. 837-842.
- [5] STUTZMAN W.L., THIELE G.A., *Antenna Theory and Design*. 3rd Edition, 2012.
- [6] RAPPAPORT T.S., *Wireless Communications: Principles & Practice*. 2nd ed, Prentice Hall PTR, Upper Saddle River, NJ, 2001.

- [7] ARCIUCH A., *Testing The Ability of Smartphones to Locate Devices With A Bluetooth Interface*. Proceedings of the 37th International Business Information Management Association (IBIMA), 30-31 May 2021, Cordoba, Spain, 2021.
- [8] HEYN R. et al., *User Tracking for Access Control with Bluetooth Low Energy*. IEEE 89th Vehicular Technology Conference (VTC2019-Spring), 2019, pp. 1-7.
- [9] HUH J.H., SEO K., *An Indoor Location-Based Control System Using Bluetooth Beacons for IoT Systems*. Sensors 2017, 17(12), 2917.
- [10] RIDA M.E., LIU F., JADI Y., ALGAWHARI A.A. ASKOURIH A., *Indoor Location Position Based on Bluetooth Signal Strength*. 2nd International Conference on Information Science and Control Engineering, 2015, pp. 769-773.

Evaluation of the accuracy and repeatability of the RSSI parameter measurement in BLE systems

ABSTRACT: Access control systems are widely used in organizations. Specialized cards and readers are used to control access. The use of cards and readers equipped with a BLE interface in such systems eliminates the need to bring the card closer to the reader in order to access the organization's resource. The paper evaluates the usability of devices with the BLE interface in automatic access control systems, based on the accuracy and repeatability of the RSSI parameter measurement.

KEYWORDS: BLE, RSSI, distance measurement

Praca wpłynęła do redakcji: 7.10.2021 r.

Recenzja książki „Komunikacja sieciowa: źródła informacji big data” autorstwa Dariusza Jarugi

Antoni DONIGIEWICZ

Instytut Teleinformatyki i Cyberbezpieczeństwa, Wydział Cybernetyki, WAT,
ul. gen. Sylwestra Kaliskiego 2, 00-908 Warszawa
antoni.donigiewicz@wat.edu.pl

STRESZCZENIE: W artykule przedstawiono recenzję książki „Komunikacja sieciowa: źródła informacji big data”. Przedstawiono rozdziały książki z krótką ich charakterystyką. Podstawowa część książki to opis systemu Jazon wykorzystanego przez autora do badań.

SŁOWA KLUCZOWE: przetwarzanie informacji big data, system Jazon

Wprowadzenie

Recenzowana książka autorstwa Dariusza Jarugi jest jedną z nielicznych (może nawet dotychczas jedyną¹) książek polskiego autora dotyczących problematyki „big data” [1]. Autor jest adiunktem na Wydziale Dziennikarstwa, Informacji i Bibliologii Uniwersytetu Warszawskiego i twórcą oprogramowania wykorzystywanego do zbierania i przetwarzania informacji big data. Sam termin „big data” rozpowszechniony przez Johna R. Mashey’a [2] w latach 90. ubiegłego wieku jest powszechnie stosowany na określenie zbiorów danych cechujących się dużą różnorodnością i docierających do odbiorców w coraz większych ilościach i z coraz większą szybkością. Zbiory te są tak obszerne, że tradycyjne oprogramowanie do przetwarzania danych po prostu nie jest w stanie nimi

¹ Poza pracą zbiorową „Big data w zarządzaniu” J. Wieczorkowski, I. Chomiak-Orsa, I. Pawełoszek, która ukazała się w tym samym roku [3].

zarządzać. Te ogromne ilości danych można wykorzystywać do rozwiązywania problemów biznesowych, z którymi wcześniej przedsiębiorstwo nie byłoby sobie w stanie poradzić. Jak się można zorientować z lektury książki, nie tylko problemy biznesowe można rozwiązywać, wykorzystując opracowane narzędzie o nazwie Jazon. System był wykorzystywany na potrzeby badań z obszaru socjologii, mediów, poligrafii, bankowości, wizerunku osób i instytucji, rynku kryptowalut i technologii blockchain.

1. Charakterystyka ogólna

Książka dotyczy problemu gromadzenia wybranych informacji ze zbioru Big Data i ich przekazywania do dalszego przetwarzania – do rafinacji informacji. Przedmiotem opracowania jest rozwiązanie problemu selekcji informacji w formie cyfrowej i ich przekazywanie do wykorzystania przez ludzi. Autor skupił się na aspektach dotyczących źródeł informacji i kolekcjonowania danych na potrzeby procesu rafinacji informacji oraz normalizacji informacji w bazie danych i udostępnianiu informacji do wykorzystania przez ludzi. Naczelnym celem książki jest pogłębiona analiza metod i narzędzi, umożliwiających realizację badań w zakresie dużych zbiorów informacji określanych jako Big Data.

Tematyka opracowania to wskazanie ogólnych ram programowych narzędzi informatycznych realizujących przetwarzanie w zakresie Big Data sprowadzające się do:

- identyfikacji źródeł informacji,
- kolekcjonowania danych (pobieranie treści ze źródeł),
- normalizacji – przekształcania różnych form informacji cyfrowej na postać znormalizowaną („oczyszczoną”),
- tworzenia bazy danych zawierającej zgromadzone dane,
- udostępnienia informacji na potrzeby procesów rafinacji danych dla systemów i narzędzi zewnętrznych.

2. Rozdział pierwszy – Big Data we współczesnym świecie

W pierwszym rozdziale autor opisuje Internet i media społecznościowe jako źródła informacji, wskazując na historyczny ich rozwój. Opis obejmuje media masowe jednokierunkowe oraz media społecznościowe dwukierunkowe, ze wskazaniem na najpopularniejsze serwisy. Przedstawiono główne funkcje serwisów społecznościowych oraz podano sposoby i możliwości pozyskiwania informacji ze stron serwisów. Autor podaje rekomendowane sposoby pobierania

treści ze stron serwisów. Niezależnie od serwisu społecznościowego, w każdym z nich występują pewne stałe elementy. Scharakteryzowano tutaj:

- profile użytkowników / konta, listy kontaktów, grupy, kręgi;
- współdzielenie plików i współdzielenie linków;
- mikroopinie, komentarze i fora dyskusyjne;
- chat i komunikatory oraz blogi i mikroblogi;
- wewnętrzne aplikacje w tym gry, quizy, ankiety itp.;
- wewnętrzną pocztę elektroniczną i powiadomienia;
- statystyki, system uprawnień;
- streaming wideo / audio;
- panel administracyjny i API (Application Programming Interface).

Autor przedstawia tutaj również możliwości informacyjnego zasilania Big Data z urządzeń Internetu rzeczy (IoT) oraz rozwija pojęcie Big Data, podając definicje dostępne w podstawowych pozycjach literatury oraz definicje proponowane przez wiodące światowe firmy, takie jak IBM, SAS Institute czy IEEE. Mamy tutaj metody pozyskiwania nowych informacji z zasobów Big Data. Analiza danych Big Data pozwala na określenie ogólnej cechy lub kierunku zjawiska, niż poznanie go w szczegółach. Pokazano, jak można identyfikować osobę na podstawie pozostawianych śladów przez urządzenia m.in. IoT.

Trzeba zdawać sobie sprawę z tego, że dane prywatne są dostępne i przetwarzane niezależnie od tego czy zostały udostępnione świadomie (wpisane do sieci), czy nieświadomie. Pokazano, jak łatwy jest dostęp do danych klienta. Wskazano na niebezpieczeństwa na styku prywatności i ochrony danych osobowych oraz nieroztropnego korzystania z sieci w warunkach bardzo dużych możliwości technologicznych gromadzenia i analizy danych. Źródłem informacji Big Data mogą być (i są najczęściej) strony internetowe o dostępie otwartym. W źródłach może być używany język polski lub angielski – nie ma to znaczenia dla technologii pozyskiwania danych. Pobieranie danych z sieci realizowane jest bez zakłócania pracy źródła danych.

3. Rozdział drugi – System Big Data Jazon

Drugi rozdział książki to opis systemu Jazon. Uzasadnienie budowy systemu to przede wszystkim bezpośredni dostęp do danych Big Data, realizowanie badań na dużym zbiorze danych, możliwość powtórzenia badań, możliwość wykorzystania narzędzia do innych badań na danych Big Data i możliwość wyznaczania prawdopodobieństwa zjawisk na podstawie zebranych danych. W tej części książki przedstawione są też formaty źródeł danych w

Internecie z krótkim opisem, przy czym sygnalizowany jest problem z plikami zabezpieczonymi hasłem i mechanizmem CAPTCHA. Autor wskazuje również na ograniczanie dostępu do informacji, czyli cenzurę, która ma istotny wpływ na możliwości kolekcjonowania danych. Podobnie jak wymagania wyrażenia zgody na przetwarzanie informacji stosowane przez niektóre źródła. Podawane są pewne możliwości rozwiązania takich problemów, ale wymagają one dodatkowych indywidualnych przedsięwzięć, które bywają dość kosztowne. Przedstawiono również występowanie zjawiska duplikatów informacji i działania eliminujące takie duplikaty.

Dla systemu Jazon poza siecią, alternatywnym źródłem danych są fizyczne nośniki danych – cyfrowe i analogowe. Część materiałów dostępnych w Internecie zawiera dane tekstowe pochodzące z materiałów drukowanych, które zostały zeskanowane i bez procesu OCR zapisane w plikach graficznych. Autor pokazuje, jakie są w tym zakresie możliwości uwzględnienia danych z takich nośników i materiałów.

System Big Data Jazon to cztery podstawowe moduły: zbierania danych, monitorowania, wykonywania kopii bezpieczeństwa i eksportu danych. Podstawowym jest moduł zbierania danych – od niego zależy w dużej mierze powodzenie procesu badawczego. Celem systemu Jazon jest odciążenie badacza od ręcznego kolekcjonowania danych. Trzeba jednak pamiętać, że na badaczu spoczywa odpowiedzialność, która zależy od doświadczenia badacza i jego intuicji w zakresie doboru wartościowych, z punktu badania, źródeł informacji.

Monitorowanie i kolekcjonowanie treści źródeł informacyjnych obejmuje przygotowanie listy adresów URL (linków) do stron WWW oraz pozostałych źródeł danych cyfrowych innych usług IT oraz zbiorów cyfrowych w postaci plików, włącznie z tymi, które są przechowywane lokalnie. Źródłami informacji są też multiwyszukiwarki, informacje patentowe, zbiory danych publicznych udostępnianych przez instytucje państwowe, różne serwisy (np. dotyczące rynku giełdowego) itp. Podstawowym komponentem modułu zbierania danych jest bot (agent). Moduł kolekcjonowania odpowiedzialny jest za trwały zapis do bazy danych pobranej przez agenta systemu Jazon informacji. Dane te są magazynowane w relacyjnej bazie danych rozproszonej na kilka serwerów.

Metodyka badań z użyciem narzędzia Jazon składa się następujących etapów:

1. Zdefiniowanie tematu badań wraz z założeniami i hipotezą.
2. Określenie metadanych, słów i wyrażeń kluczowych – słupów.
3. Określenie źródeł danych cyfrowych i opcjonalnie analogowych.
4. Digitalizacja/konwersja danych analogowych (opcjonalnie).
5. Opracowanie dedykowanego modułu importu dla systemu (opcjonalnie).
6. Zebranie danych cyfrowych i ich zapisanie do bazy danych systemu.

7. Opracowanie dedykowanego modułu rafinacji danych (opcjonalnie).
8. Wstępna rafinacja (opcjonalnie).
9. Rafinacja.
10. Decyzja o wykorzystaniu wyników rafinacji.
11. Analiza wyników, weryfikacja hipotezy i wniosków.

Autor szczegółowo pokazuje metodykę badań z użyciem narzędzia Jazon, opisując kolejne etapy i realizowane procedury. Podaje również konkretne przyjęte rozwiązania dla przykładowych badań.

4. Rozdział trzeci – System Jazon. Badania

W tej części autor przedstawia cztery badania zrealizowane z wykorzystaniem informacji zgromadzonych za pomocą systemu Jazon.

Pierwsze badanie dotyczyło predykcji wyników wyborów prezydenckich i parlamentarnych z 2015 roku w Polsce. Przedmiotem analizy były treści (ponad 2 miliony wpisów) zebrane od 16 maja do października 2015 roku. Analizowano wpisy pozytywne i negatywne na dany temat, wykorzystując różne źródła informacji. W opisie badań omówiono założenia, sposoby wyboru nazw i sentymentów, tryb zbierania wpisów i wyniki analizy sentymentów. W wyniku analizy wykazano istotny związek pomiędzy liczbą pozytywnych wpisów na blogach a liczbą głosów uzyskanych w wyborach parlamentarnych. Uzyskane w badaniu wyniki są niemal identyczne z tymi, które otrzymano podczas tradycyjnych badań sondażowych i z wynikami głosowania ogłoszonymi przez Państwową Komisję Wyborczą (wystąpiły niewielkie różnice poniżej 1%).

Drugie z przedstawionych badań dotyczyło kolekcjonowania informacji historycznych, niezbędnych do przeprowadzenia badania autorytetu Jana Pawła II. Przedstawiono sposoby ustalenia różnych historycznych źródeł informacji i sposoby wyodrębniania daty publikacji informacji w źródle. Wykorzystano publikacje z lat 1996-2015. Przedmiotem szczegółowych analiz i wniosków były zebrane przez robota Big Data dane z ponad dziesięcioletniego okresu. Źródłami były grupy dyskusyjne USENET i serwisy WWW – zebrano ponad 4,7 miliona wpisów. Zebrane dane były opracowywane przez specjalistę z zakresu nauk społecznych. Do dalszej analizy wyodrębniono w ramach procesu rafinacji sieciowej 55 747 wpisów – wiadomości odpowiadających przyjętym przez specjalistę założeniom. Dużą trudnością było powiązanie wiadomości z datą jej wystąpienia. Zebrane dane pozwoliły na określenie, jak zmieniały się wraz z upływem czasu opinie na temat Jana Pawła II.

Pokazano w ten sposób, że system Jazon może być użytecznym narzędziem do badania przeszłości, pod warunkiem, że dane z przeszłego okresu są dostępne online lub w postaci zdigitalizowanych dokumentów.

Trzecie z opisywanych badań dotyczyło identyfikacji przyczyn chorób Polaków. Bazą rafinacji były informacje zgromadzone za pomocą systemu Jazon. Przed kolekcjonowaniem wpisów ze stron za pomocą Jazona, bazę monitorowanych źródeł uzupełniono o dedykowane adresy sieciowe otwartych źródeł danych, które dotyczyły przedmiotowego badania – zdrowia. Do badania, poza źródłami informacji, które system Jazon kolekcjonował wcześniej, dodano nowe, bezpośrednio związane z tematyką badania, dotyczące zdrowia i chorób. W opisie przedstawiono szczegóły dotyczące badania, w tym opracowania zbioru nazw i sentymentów, kolekcjonowania wpisów ze stron, weryfikacji sentymentów, obliczania frekwencji i interpretacji wyników. Otrzymane szczegółowe wyniki rafinacji sieciowej pokrywają się z raportem WHO z 2002 roku dotyczącym determinantów mających bezpośredni wpływ na zdrowie człowieka. Pokazano, że nawet krótkie, bo dwutygodniowe zbieranie danych – przy wykorzystaniu poprawnej metodyki i dostosowanych do niej narzędzi może przynieść wyniki, które pokrywają się z danymi zawartymi w tradycyjnych, autorytatywnych wynikach badań.

Czwarte z przedstawionych badań dotyczyło wykrycia nieuczciwych praktyk stosowanych przez niektórych uczestników akcji „Wsparcie dla bibliotek”. Opisywane badanie dotyczyło danych ustrukturyzowanych, to jest takich, które mają jednoznacznie określoną budowę. Sprawdzano, czy analiza danych ze stron konkursowych, zebranych okresowo przez system Jazon, bazując wyłącznie na ogólnie dostępnych danych, pozwoli wykryć nieuczciwe praktyki uczestników akcji. Monitorowano oddawane głosy na szkoły biorące udział w konkursie. Przedstawiono procedury zastosowane w celu zebrania przez robota danych ustrukturyzowanych i ich analizy. Dane były pobierane przez system Jazon z różną częstotliwością uzależnioną od dynamiki zmian liczby oddawanych głosów. Zebrane przez robota dane były szczytkowe, w porównaniu z tym, czym dysponował organizator, a mimo to pozwoliły zarejestrować anomalie, które znalazły w późniejszym czasie swoje potwierdzenie w wynikach ogłoszonych oficjalnie przez organizatora konkursu. Przeprowadzona analiza pokazała nową drogę badań, w których możliwe jest wyciąganie wniosków na podstawie zbieranych w czasie rzeczywistym danych ustrukturyzowanych np. z urzędów pomiarowych, kursów, notowań itp. Należy jednak pamiętać, że dla każdego badania, w którym dane mają postać ustrukturyzowaną, trzeba napisać program, który we właściwy sposób wyodrębni informacje niezbędne do dalszego badania.

Podsumowanie

W aneksach książki przedstawiono problemy przesyłania dużych ilości danych pomiędzy elementami systemu oraz podano charakterystykę techniczną modułów systemu Jazon i problemy związane z zapewnieniem wydajności, zarządzania, rozwoju systemu i bezpieczeństwa przechowywanych danych.

Bibliografia opracowania, bardzo obszerna i wskazująca na dobrą znajomość tematyki przez autora, została podzielona na trzy części: wydawnictwa zwarte, artykuły i netografię. Wydawnictwa zwarte to dwadzieścia osiem pozycji głównie w języku polskim. Artykuły to trzydzieści siedem pozycji. Netografia natomiast jest bardzo obszerna, bo obejmuje 304 pozycje w znacznej części w języku angielskim. Na te pozycje głównie powołuje się autor przy analizie tematyki Big Data w bardzo obszernym pierwszym rozdziale. Poza wskazanymi pozycjami bibliografii w załączniku zamieszczono literaturę przedmiotu, którą wykorzystano podczas budowy robota Jazon. Literatura to łącznie 77 pozycji w znacznej części w języku angielskim, dotyczy źródeł internetowych i dokumentacji systemów oraz programów użytych w konstrukcji robota Jazon.

W książce pokazano metodyki kolekcjonowania użytecznych informacji na potrzeby procesu rafinacji. Szczególny nacisk położono na źródła informacji oraz ich kolekcjonowanie. Pokazano kolekcjonowanie danych z dużych nieustrukturyzowanych i ustrukturyzowanych źródeł informacji. Wyniki przedstawionych badań potwierdziły przydatność opracowanej metody i narzędzia informatycznego do gromadzenia informacji źródłowych oraz potwierdziły, że system Jazon wraz z metodyką rafinacji danych jest skutecznym narzędziem badawczym.

Książka będzie pomocą dla wszystkich zainteresowanych badaniami w obszarze Big Data, w szczególności dla osób rozpoczynających realizację badań oraz studentów realizujących prace dyplomowe z tego zakresu.

Literatura

- [1] JARUGA D., *Komunikacja sieciowa: źródła informacji big data*. Wydawnictwo Naukowe i Edukacyjne Stowarzyszenia Bibliotekarzy Polskich, Warszawa, 2021.
- [2] MASHEY R.J., *Big Data... and the Next Wave of InfraStress*. Slides from invited talk. Usenix. Retrieved September 28, 2016.
- [3] WIECZORKOWSKI J, CHOMIAK-ORSA I., PAWEŁOSZEK I., *Big data w zarządzaniu*. Wydawnictwo PWE, Warszawa, 2021.

Book review by Dariusz Jaruga

“Network communications: Big Data information sources”

ABSTRACT: This paper reviews book “Network Communications: Big Data information sources”. Chapters of the book are presented along with a brief characteristic of them. The core part of the book is a description of the Jazon system, used by the author for research.

KEYWORDS: big data information processing, Jazon system

Praca wpłynęła do redakcji: 30.12.2021 r.

Information for Authors
– rules of papers preparation and reviewing for
TELEINFORMATICS REVIEW

The *Teleinformatics Review* is devoted to the publication of original research results in fields of science including, but not limited to: computer science, telecommunication, signal processing, network systems, automation and robotics, etc., which have not been published elsewhere in their entirety or considerable part. If a submitted paper is a part of another published work, e.g. a doctoral dissertation, a postdoctoral thesis, etc., the source work should be included in the list of literature and the editorial office must be informed about it.

In order to publish a paper in the Teleinformatics Review it is necessary to submit it to the editorial office in an electronic form (and possibly its printed copy, one-sided, legible, on white A4 sheets) according to the given template. Only original works in English or Polish will be accepted. The text of the paper should be prepared in the format of Microsoft Word editor (versions 2003 or 2010 are suggested). Appropriate templates can be downloaded from website review.ita.wat.edu.pl (or przeglad.ita.wat.edu.pl). The electronic version submitted to the editorial office should contain a source file of the paper in DOC or DOCX format, with all figures and tables being inserted. The editorial office does not rewrite the text neither make drawings. In addition to the mentioned source file, all figures should be delivered in commonly used image formats (preferably as EPS, JPG, TIFF, or others).

Papers to be published in the Teleinformatics Review are subject to initial acceptance by the editorial office and then are subject to review by two external reviewers. Reviewers and authors do not know each other personal data. The content of the review will be available at the editorial office. If one review is negative (or imprecise) then a third reviewer may be appointed. If both reviews are negative the paper is rejected. If the review indicates a necessity of some corrections, the author must consider all of them and resubmit the improved paper by the determined deadline.

The volume of a submitted paper generally not exceed 20 pages of typescript A4. A deviation from this rule requires agreement of the editorial office. Except the last page, no more than 10% of any page within the paper can be left empty. Figures must be numbered and described below them as well as tables must be numbered and described at the top of them. The literature should hold the form given in the template.

The authors are obliged to submit a statement to the editorial office on the percentage contribution to the creation of the accepted paper, confirming the lack

of prior publication of such a work, or a public speech on the subject at a conference or symposium.

The editorial board reserves rights to introduce minor editorial changes to the content of paper without consulting the author. The editorial office insists that no special formatting should be used, which would be inconsistent with the template.

Papers printed in the Teleinformatics Review and their abstracts are placed in the national database of Polish technical journals BazTech as well as on the INDEX COPERNICUS website. Additionally, the papers will be available in the electronic PDF form on website review.ita.wat.edu.pl.

Publication in the Teleinformatics Review does not involve any costs for authors. The editorial office does not charge for submitting, reviewing, preparing for publication and publishing the work. The publication of a paper in the Teleinformatics Review is tantamount to transfer of authors' property rights for publication to the publisher, i.e. the Military University of Technology. By submitting a paper for publication in the Teleinformatics Review, the author agrees, for publication purposes, to the processing by the editorial office the author's name, email address, affiliation, and other contact details.



All papers published in the journal **TELEINFORMATICS REVIEW** are made available under the Creative Commons Attribution-NonCommercial-NoDerivatives 3.0 (CC BY-NC-ND 3.0) license. Thus, licensees may copy, distribute, display, and perform the work and make derivative works and remixes based on it only for non-commercial purposes; licensees may copy, distribute, display and perform only verbatim copies of the work, not derivative works and remixes based on it.

The editorial office does not return received materials.

The editorial office does not pay fees for papers publishing.

The editor-in-chief may refuse to publish a paper in the following cases:

- if the content of the paper violates the law (principles of secrecy protection, press law, copyright law, etc.) or good manners;
- the author does not agree to introduce all necessary corrections proposed by the editorial board or reviewers;
- the text and illustrative material submitted by the author does not meet the technical requirements given in this document or the template.

Informacje dla autorów
– zasady przygotowania tekstu i recenzowania artykułów do
PRZEGLĄDU TELEINFORMATYCZNEGO

W Przeglądzie Teleinformatycznym zamieszczane są oryginalne artykuły z dziedzin: *informatyka, telekomunikacja, przetwarzanie sygnałów, systemy sieciowe, automatyka i robotyka* oraz pokrewnych, niepublikowane dotychczas w całości lub w znaczącej części. Jeśli nadesłana praca stanowi część innej opublikowanej pracy, np. pracy doktorskiej, habilitacji, etc., to źródło powinno być umieszczone w spisie literatury, a redakcja powinna być o tym poinformowana.

W celu opublikowania artykułu w *Przeglądzie* niezbędne jest dostarczenie do redakcji treści artykułu w postaci **elektronicznej** według podanego szablonu i ewentualnie jednego egzemplarza wydrukowanego (jednostronnie, czytelnie, na białym papierze formatu A4). Przyjmowane są tylko oryginalne prace w języku angielskim lub polskim. Tekst artykułu powinien być przygotowany w formacie edytora Microsoft Word (wersja 2003 lub 2010 jest zalecana). Szablony dla artykułów są dostępne w pliku na stronie przeklad.ita.wat.edu.pl (lub review.ita.wat.edu.pl). Przekazane do redakcji materiały powinny zawierać plik źródłowy w formacie DOC lub DOCX, z wstawionymi rysunkami. Redakcja nie przepisuje tekstów i nie wykonuje rysunków. Dodatkowo należy dostarczyć pliki źródłowe rysunków (najlepiej w formacie EPS, JPG, TIFF lub innym powszechnie używanym).

Artykuły przeznaczone do opublikowania w *Przeglądzie* podlegają wstępnej ocenie przez redaktora działu, a następnie podlegają recenzji przez dwóch zewnętrznych recenzentów. Recenzenci i autorzy nie znają swoich danych personalnych. Z treścią recenzji można zapoznać się w redakcji. Jeśli jedna z recenzji jest negatywna (lub nieprecyzyjna), może być powołany trzeci recenzent. Jeśli dwie recenzje są negatywne, artykuł jest odrzucany. Jeśli z recenzji wynika konieczność dokonania poprawek w treści artykułu, to autor jest zobowiązany do ich rozpatrzenia i dostarczenia do redakcji poprawionej wersji artykułu, w terminie ustalonym przez redakcję.

Objętość artykułu zasadniczo nie powinna przekroczyć 20 stron maszynopisu A4. Odstąpienie od tej zasady wymaga uzgodnień z redakcją *Przeglądu*. Na stronach tekstu artykułu nie może być pozostawione więcej niż 10% pustego miejsca, za wyjątkiem ostatniej strony. Rysunki należy numerować i opatrzyć (pod spodem) wyczerpującym podpisem. Tabele również muszą być numerowane (tytuł nad tabelą). Literatura może być uszeregowana alfabetycznie oraz powinna mieć postać jak w szablonie.

Autorzy są zobligowani do złożenia w redakcji oświadczenia autorskiego o wkładzie procentowym w powstanie artykułu, braku wcześniejszej publikacji artykułu w przedstawionej formie lub wystąpieniu publicznym na ten temat na konferencji lub sympozjum.

Redakcja zastrzega sobie prawo wprowadzenia niewielkich redakcyjnych zmian w treści artykułu bez konsultacji z autorem. Redakcja nalega, aby **nie stosować** żadnego specjalnego formatowania i **trzymać się ściśle** ustaleń zawartych w szablonie.

Streszczenia i pełne teksty artykułów drukowanych w *Przeglądzie* zamieszczane są w krajowej bazie danych o zawartości polskich czasopism technicznych BazTech oraz na platformie INDEX COPERNICUS. Opublikowane w *Przeglądzie* artykuły będą także w całości udostępnione w internetowej wersji (format PDF) czasopisma, pod adresem przeglad.ita.wat.edu.pl (lub review.ita.wat.edu.pl).

Publikacja w *Przeglądzie* nie wiąże się z żadnymi kosztami dla autorów. Redakcja nie pobiera opłat za zgłoszenie, przygotowanie do druku, recenzję czy publikację pracy. Przekazanie artykułu do publikacji w *Przeglądzie* jest równoznaczne z przekazaniem autorskich praw majątkowych do publikacji na rzecz wydawcy, tj. Wojskowej Akademii Technicznej. Przekazując artykuł do publikacji w *Przeglądzie* autor zgadza się na przechowywanie i przetwarzanie przez redakcję, w celach publikacyjnych, imienia, nazwiska, adresu e-mail i afiliacji.



Wszystkie artykuły opublikowane w czasopiśmie **PRZEGLĄD TELEINFORMATYCZNY (TELEINFORMATICS REVIEW)** są udostępniane na licencji Creative Commons Uznanie autorstwa – Użycie niekomercyjne – Bez utworów zależnych 3.0 (CC BY-NC-ND 3.0), która zezwala na kopiowanie, przedstawianie i rozpowszechnianie utworu jedynie w celach niekomercyjnych oraz pod warunkiem zachowania go w oryginalnej postaci (czyli nietworzenia utworów zależnych), przy jednoczesnym odpowiednim oznaczeniu autorstwa utworu.

Redakcja nie zwraca materiałów dostarczonych do redakcji.

Redakcja nie przewiduje honorariów za opublikowanie artykułu.

Redaktor naczelny może odmówić opublikowania artykułu w przypadku, gdy:

- treści zawarte w materiałach naruszają prawo (zasady ochrony tajemnicy, prawo prasowe, prawo autorskie itp.) lub dobre obyczaje;
- autor nie zgadza się na wprowadzenie wszystkich koniecznych poprawek zaproponowanych przez redakcję lub recenzentów;
- tekst i materiał ilustracyjny złożony przez autora nie spełnia wymagań technicznych podanych w niniejszym dokumencie i szablonie.