

# PRACE NAUKOWE

Uniwersytetu Ekonomicznego we Wrocławiu

# RESEARCH PAPERS

of Wrocław University of Economics

Nr 328

**Taksonomia 23**

**Klasyfikacja i analiza danych –  
teoria i zastosowania**

Redaktorzy naukowci

Krzysztof Jajuga, Marek Walesiak



Wydawnictwo Uniwersytetu Ekonomicznego we Wrocławiu  
Wrocław 2014

Redaktor Wydawnictwa: Barbara Majewska

Redaktor techniczny: Barbara Łopusiewicz

Korektor: Barbara Cibis

Łamanie: Beata Mazur

Projekt okładki: Beata Dębska

Publikacja jest dostępna w Internecie na stronach:

[www.ibuk.pl](http://www.ibuk.pl), [www.ebscohost.com](http://www.ebscohost.com),

w Dolnośląskiej Bibliotece Cyfrowej [www.dbc.wroc.pl](http://www.dbc.wroc.pl),

The Central and Eastern European Online Library [www.ceeol.com](http://www.ceeol.com),

a także w adnotowanej bibliografii zagadnień ekonomicznych BazEkon

[http://kangur.uek.krakow.pl/bazy\\_ae/bazekon/nowy/index.php](http://kangur.uek.krakow.pl/bazy_ae/bazekon/nowy/index.php)

Informacje o naborze artykułów i zasadach recenzowania znajdują się  
na stronie internetowej Wydawnictwa

[www.wydawnictwo.ue.wroc.pl](http://www.wydawnictwo.ue.wroc.pl)

Tytuł dofinansowany ze środków Narodowego Banku Polskiego  
oraz ze środków Sekcji Klasyfikacji i Analizy Danych PTS

Kopiowanie i powielanie w jakiegokolwiek formie  
wymaga pisemnej zgody Wydawcy

© Copyright by Uniwersytet Ekonomiczny we Wrocławiu  
Wrocław 2014

**ISSN 1899-3192** (Prace Naukowe Uniwersytetu Ekonomicznego we Wrocławiu)

**ISSN 1505-9332** (Taksonomia)

Wersja pierwotna: publikacja drukowana

Druk: Drukarnia TOTEM

## Spis treści

<b>Wstęp</b> .....	11
<b>Małgorzata Rószkiewicz</b> , Wykorzystanie metaanalizy w budowaniu modelu pomiarowego w przypadku braku niezmienniczości zasad pomiaru na przykładzie pomiaru zadowolenia z życia.....	13
<b>Elżbieta Sobczak</b> , Harmonijność inteligentnego rozwoju regionów Unii Europejskiej .....	21
<b>Ewa Roszkowska, Renata Karwowska</b> , Analiza porównawcza województw Polski ze względu na poziom zrównoważonego rozwoju w roku 2010.....	30
<b>Tadeusz Kufel, Magdalena Osińska, Marcin Błażejowski, Paweł Kufel</b> , Analiza porównawcza wybranych filtrów w analizie synchronizacji cyklu koniunkturalnego.....	41
<b>Marcin Salamaga</b> , Próba konstrukcji tablic „wymierania scenicznego” spektakli operowych na przykładzie Metropolitan Opera.....	51
<b>Iwona Foryś</b> , Wykorzystanie analizy dyskryminacyjnej do typowania rynków podobnych w procesie wyceny nieruchomości niemieszkalnych .....	59
<b>Jerzy Korzeniewski</b> , Selekcja zmiennych w klasyfikacji – propozycja algorytmu .....	69
<b>Sabina Denkowska</b> , Testowanie wielokrotne przy weryfikacji wieloczynnikowych modeli proporcjonalnego hazardu Coxa.....	76
<b>Ewa Chodakowska</b> , Teoria równań strukturalnych w klasyfikacji zmiennych jawnych i ukrytych według charakteru ich wzajemnych oddziaływań .....	85
<b>Iwona Konarzewska</b> , Model PCA dla rynku akcji – studium przypadku .....	94
<b>Katarzyna Wójcik, Janusz Tuchowski</b> , Dobór optymalnego zestawu słów istotnych w opiniach konsumentów na potrzeby ich automatycznej analizy	106
<b>Aleksandra Łuczak</b> , Zastosowanie metody AHP-LP do oceny ważności determinant rozwoju społeczno-gospodarczego w jednostkach administracyjnych .....	116
<b>Aleksandra Witkowska, Marek Witkowski</b> , Klasyfikacja pozycyjna banków spółdzielczych według stanu ich kondycji finansowej w ujęciu dynamicznym .....	126
<b>Adam Depta</b> , Zastosowanie analizy korespondencji do oceny jakości życia ludności na podstawie kwestionariusza SF-36v2 .....	135
<b>Marek Lubicz, Maciej Zięba, Konrad Pawelczyk, Adam Rzechonek, Marek Marciniak, Jerzy Kołodziej</b> , Indukcja reguł dla danych niekompletnych i niezbalansowanych: modele klasyfikatorów i próba ich zastosowania do predykcji ryzyka operacyjnego w torakochirurgii .....	146

<b>Małgorzata Misztal</b> , Wybrane metody oceny jakości klasyfikatorów – przegląd i przykłady zastosowań.....	156
<b>Anna M. Olszewska</b> , Wykorzystanie wybranych metod taksonomicznych do oceny potencjału innowacyjnego województw .....	167
<b>Iwona Bąk</b> , Porównanie jakości grupowań powiatów województwa zachodniopomorskiego pod względem atrakcyjności turystycznej.....	177
<b>Agnieszka Kozera, Joanna Stanisławska, Romana Głowicka-Wołoszyn</b> , Segmentacja gospodarstw domowych według wydatków na turystykę zorganizowaną.....	186
<b>Agnieszka Wałęga</b> , Podejście syntetyczne w analizie spójności ekonomicznej gospodarstw domowych.....	196
<b>Joanna Banaś, Małgorzata Machowska-Szewczyk, Bożena Mroczek</b> , Zastosowanie analizy korespondencji do badania wpływu elektrowni wiatrowych na jakość życia ludności .....	205
<b>Joanna Banaś, Krzysztof Małecki</b> , Klasyfikacja punktów pomiarów ankietowych kierowców na granicy Szczecina z wykorzystaniem zmiennych symbolicznych.....	214
<b>Aneta Becker</b> , Wykorzystanie informacji granularnej w analizie wymagań rynku pracy.....	222
<b>Katarzyna Cheba, Joanna Holub-Iwan</b> , Wykorzystanie analizy korespondencji w segmentacji rynku usług medycznych.....	230
<b>Adam Depta, Iwona Staniec</b> , Identyfikacja czynników decydujących o jakości życia studentów łódzkich uczelni.....	238
<b>Katarzyna Dębowska, Jarosław Kilon</b> , Reguły asocjacyjne w analizie wyników badań metodą Delphi.....	247
<b>Anna Domagała</b> , O wykorzystaniu analizy głównych składowych w metodzie <i>Data Envelopment Analysis</i> .....	254
<b>Alicja Grześkowiak</b> , Analiza wykluczenia cyfrowego w Polsce w ujęciu indywidualnym i regionalnym.....	264
<b>Anna M. Olszewska, Anna Gryko-Nikitin</b> , Pomiar postrzegania jakości kształcenia uczelni wyższej na danych porządkowych z wykorzystaniem środowiska R.....	273
<b>Karolina Paradysz</b> , Hierarchiczna metoda grupowania powiatów jako podejście benchmarkowe w ocenie bezrobocia według BAEL-u w wybranych typach małych obszarów .....	282
<b>Radosław Pietrzyk</b> , Porównanie metod pomiaru efektywności zarządzania portfelami funduszy inwestycyjnych.....	290
<b>Agnieszka Przedborska, Małgorzata Misztal</b> , Wybrane metody statystyki wielowymiarowej w ocenie skuteczności terapeutycznej głębokiej stymulacji elektromagnetycznej u pacjentów z chorobą zwyrodnieniową stawów.....	299

<b>Wojciech Roszka, Marcin Szymkowiak</b> , Podejście kalibracyjne w statystycznej integracji danych .....	308
<b>Iwona Skrodzka</b> , Zastosowanie wybranych metod klasyfikacji do analizy kapitału ludzkiego krajów Unii Europejskiej .....	316
<b>Agnieszka Stanimir</b> , Wielowymiarowa analiza czynników sprzyjających włączeniu społecznemu .....	326
<b>Dorota Strózik, Tomasz Strózik</b> , Przestrzenne zróżnicowanie poziomu życia w województwie wielkopolskim.....	334
<b>Izabela Szamrej-Baran</b> , Identyfikacja przyczyn ubóstwa energetycznego w Polsce przy wykorzystaniu modelowania miękkiego.....	343
<b>Janusz Tuchowski, Katarzyna Wójcik</b> , Klasyfikacja obiektów w systemie Krajowych Ram Kwalifikacji opisanych za pomocą ontologii .....	353
<b>Aleksandra Matuszewska-Janica</b> , Grupowanie krajów Unii Europejskiej ze względu na poziom feminizacji sektorów gospodarczych .....	361
<b>Monika Rozkrut, Dominik Rozkrut</b> , Identyfikacja strategii innowacyjnych przedsiębiorstw usługowych w Polsce .....	369

## Summaries

<b>Małgorzata Rószkiewicz</b> , The use of meta-analysis in building the measurement model in case of the absence of measurement invariance on the example of measuring of life satisfaction.....	20
<b>Elżbieta Sobczak</b> , Harmonious smart growth of European Union regions.....	29
<b>Ewa Roszkowska, Renata Karwowska</b> , The comparative analysis of Polish voivodeships with respect to sustainable development in 2010.....	40
<b>Tadeusz Kufel, Magdalena Osińska, Marcin Błażejowski, Paweł Kufel</b> , Comparative analysis of chosen filters in business cycles analysis .....	50
<b>Marcin Salamaga</b> , The attempt of construction of the life tables for opera works on the example of the Metropolitan Opera .....	58
<b>Iwona Foryś</b> , Using discriminant analysis to select similar markets in non-residential property valuation process.....	68
<b>Jerzy Korzeniewski</b> , Variable selection in classification – algorithm proposal .....	75
<b>Sabina Denkowska</b> , Multiple testing in the verification process of multifactorial Cox proportional hazards models .....	84
<b>Ewa Chodakowska</b> , The theory of structural equations modelling in the classification of observed variables and latent constructs according to the character of their relationship.....	93
<b>Iwona Konarzewska</b> , Modelling stock market by PCA factor model – case study .....	105

<b>Katarzyna Wójcik, Janusz Tuchowski</b> , Selection of the optimal set of relevant words in consumers opinions in the context of the opinion mining ..	115
<b>Aleksandra Łuczak</b> , Application of AHP-LP to the evaluation of importance of determinants of socio-economic development in the administrative units .....	125
<b>Aleksandra Witkowska, Marek Witkowski</b> , A dynamic approach to the ranking of cooperative banks by their financial condition .....	134
<b>Adam Depta</b> , Application of correspondence analysis for the measurement of quality of life – questionnaire SF-36v2 based research .....	145
<b>Marek Lubicz, Maciej Zięba, Konrad Pawelczyk, Adam Rzechonek, Marek Marciniak, Jerzy Kołodziej</b> , Classification rules extraction for missing and imbalance data: models of classifiers and initial results in the rules-based thoracic surgery risk prediction.....	155
<b>Małgorzata Misztal</b> , Selected methods for assessing the performance of classifiers – an overview and examples of applications.....	166
<b>Anna M. Olszewska</b> , The application of selected quantitative methods to the evaluation of voivodeship innovation level potential.....	176
<b>Iwona Bąk</b> , The comparison of the quality of groupings of poviats of West Pomeranian Voivodeship in terms of tourism attractiveness .....	185
<b>Agnieszka Kozera, Joanna Stanisławska, Romana Głowicka-Wołoszyn</b> , Household segmentation with respect to the expenditure on organized tourism.....	195
<b>Agnieszka Wałęga</b> , Synthetic approach in the analysis of economic coherence of households .....	204
<b>Joanna Banaś, Małgorzata Machowska-Szewczyk, Bożena Mroczek</b> , Using the correspondence analysis to examine the impact of wind turbines on the quality of life.....	213
<b>Joanna Banaś, Krzysztof Małecki</b> , Classification of measurement survey points of drivers on the boundary of Szczecin using symbolic variables...	221
<b>Aneta Becker</b> , The use granular information in the analysis of the requirements of the labor market.....	229
<b>Katarzyna Cheba, Joanna Hołub-Iwan</b> , The application of the correspondence analysis of patients segmentation on the medical service market .....	237
<b>Adam Depta, Iwona Staniec</b> , Identification of the factors that determine the quality of students life at universities in Lodz.....	246
<b>Katarzyna Dębkowska, Jarosław Kilon</b> , Association rules in the analysis of research results the Delphi method .....	253
<b>Anna Domagała</b> , About using Principal Component Analysis in Data Envelopment Analysis .....	263
<b>Alicja Grześkowiak</b> , Analysis of the digital divide in Poland at the individual and regional level .....	272

<b>Anna M. Olszewska, Anna Gryko-Nikitin</b> , Assessment of perception of quality of teaching at an institution of higher learning based on the ordinal data with the utilization of R environment.....	281
<b>Karolina Paradysz</b> , The hierarchical method of grouping poviats as a benchmark approach in the assessment of unemployment by BAEL in selected types of small areas .....	289
<b>Radosław Pietrzyk</b> , Comparison of methods of measuring the performance of investment funds portfolios.....	298
<b>Agnieszka Przedborska, Małgorzata Misztal</b> , Selected multivariate statistical analysis methods in the evaluation of efficacy of deep electromagnetic stimulation in patients with degenerative joint disease .....	307
<b>Wojciech Roszka, Marcin Szymkowiak</b> , A calibration approach in statistical data integration .....	315
<b>Iwona Skrodzka</b> , Application of some methods of classification to the analysis of human capital in the European Union.....	325
<b>Agnieszka Stanimir</b> , Multivariate analysis of social inclusion factors.....	333
<b>Dorota Strózik, Tomasz Strózik</b> , Spatial differentiation of the standard of living in Great Poland Voivodeship .....	342
<b>Izabela Szamrej-Baran</b> , Identification of fuel poverty causes in Poland using soft modelling .....	352
<b>Janusz Tuchowski, Katarzyna Wójcik</b> , Classification of objects in the National Classification Framework described by the ontology.....	360
<b>Aleksandra Matuszewska-Janica</b> , Clustering of European Union states taking into consideration the levels of feminization of economic sectors..	368
<b>Monika Rozkrut, Dominik Rozkrut</b> , Identification of service sector innovation strategies in Poland.....	379

**Katarzyna Wójcik, Janusz Tuchowski**

Uniwersytet Ekonomiczny w Krakowie

---

## **DOBÓR OPTYMALNEGO ZESTAWU SŁÓW ISTOTNYCH W OPINIACH KONSUMENTÓW NA POTRZEBY ICH AUTOMATYCZNEJ ANALIZY**

---

**Streszczenie:** Analiza opinii konsumenckich jest obszarem badań, który może mieć znaczący wpływ na rozwój działalności biznesowej. Wiele podejść do automatycznej analizy opinii konsumenckich opartych jest na podobieństwie lub też odległości między dokumentami tekstowymi. Powstaje więc problem, która z miar podobieństwa dokumentów tekstowych (lub ich odległości) jest najlepsza w przypadku charakterystycznego rodzaju tekstów, jakimi są opinie konsumentów. Głównym celem pracy jest przeprowadzenie analizy doboru optymalnego zestawu słów istotnych podczas obliczania podobieństwa (odległości) dokumentów tekstowych stosowanych na potrzeby automatycznej analizy opinii konsumenckich. Do obliczeń wykorzystane zostały: język R, program RapidMiner oraz arkusz kalkulacyjny.

**Słowa kluczowe:** *text mining*, *Web mining*, taksonomia, klasyfikacja dokumentów tekstowych, ocena jakości miar.

### **1. Wstęp**

Analiza opinii konsumenckich jest obszarem badań, który może mieć znaczący wpływ na rozwój działalności biznesowej. Znaczna liczba konsumentów przed dokonaniem wyboru towaru lub usługi przeszukuje Internet, zaznajamiając się z opiniami innych użytkowników sieci. Znalezione rekomendacje często odgrywają decydującą rolę podczas podejmowania decyzji. Z kolei po dokonaniu zakupu lub skorzystaniu z usługi użytkownik jest proszony o wystawienie opinii lub też dokonuje tego z własnej inicjatywy. Narastająca liczba opinii dostępnych w sieci wytworzyła potrzebę ich automatycznej analizy i przetwarzania.

Wiele podejść do automatycznej analizy opinii konsumenckich opiera się na analizie podobieństwa lub też odległości pomiędzy dokumentami tekstowymi. Powstaje zatem problem, która z miar podobieństwa dokumentów tekstowych (lub ich odległości) jest najlepsza w przypadku charakterystycznego rodzaju tekstów, jakimi są opinie konsumentów.



Głównym celem pracy jest przeprowadzenie analizy doboru optymalnego zestawu słów istotnych podczas obliczania podobieństwa (lub odległości) dokumentów tekstowych, stosowanego na potrzeby automatycznej analizy opinii konsumenckich. Działanie takie pozwoli na zwiększenie efektywności badań poprzez zawężenie analizowanych danych. Wylimitowane z analizy zostaną słowa nieistotne, niemające wpływu na wynik końcowy.

## 2. Automatyczna analiza opinii konsumenckich i miary podobieństwa tekstów

Automatyczna analiza opinii konsumenckich (*sentiment analysis, opinion mining*) to ogół działań mających na celu zautomatyzowanie procesu wyszukiwania, ekstrakcji i analizy danych pochodzących ze specyficznych tekstów, jakimi są opinie użytkowników. Są to działania z pogranicza przetwarzania języka naturalnego (*Natural Language Processing – NLP*), lingwistyki komputerowej (*computational linguistics*) oraz eksploracyjnej analizy tekstu (*text mining*). Jej celem jest określenie nastawienia autora wypowiedzi do jej przedmiotu [Wikipedia].

W ramach automatycznej analizy opinii konsumenckich wyróżnić można trzy rodzaje działań [Liu 2007; Pang, Lee 2008]:

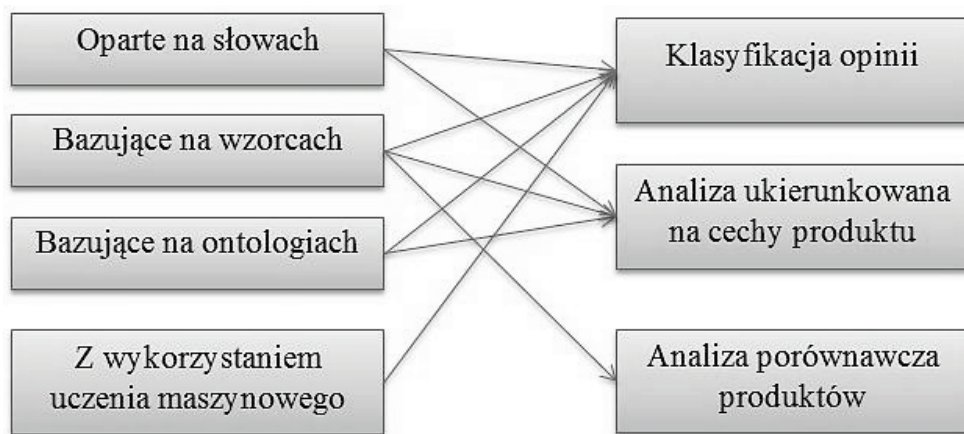
- **Klasyfikacja opinii** – podział opinii na grupy według ich nacechowania (np. pozytywne, negatywne, neutralne) lub przypisanie pojedynczej opinii jej polaryzacji (przydzielenie jej do jednej z uprzednio wymienionych grup). Brana jest tu pod uwagę opinia jako całość.
- **Analiza ukierunkowana na cechy produktu** – wyszukanie w opinii poszczególnych aspektów (cech) przedmiotu opinii, a następnie zbadanie stosunku autora wypowiedzi do tego właśnie aspektu. Badana jest nie cała opinia, ale poszczególne jej części odnoszące się do kolejnych cech opisywanego produktu czy usługi.
- **Analiza porównawcza produktów** – badanie opinii na temat jednego produktu, określonej poprzez analizę zdania porównującego go do innego produktu. Konieczne jest zidentyfikowanie w opinii zdań porównujących, a następnie ich analiza ukierunkowana na przedmiot porównania.

Rysunek 1 przedstawia różne text miningowe podejścia do automatycznej analizy opinii konsumenckich. Są one przyporządkowane do poszczególnych działań wyróżnionych w ramach automatycznej analizy opinii konsumenckich.

Koncentrując się na klasyfikacji opinii, można zauważyć, że wszystkie cztery text miningowe podejścia do automatycznej analizy opinii konsumentów znajdują w niej zastosowanie [Lula, Wójcik 2011]:

- **Podejście oparte na słowach (*word-based approach*)** – podstawą tego podejścia jest przekonanie, że znaczenie wypowiedzi (również jej nacechowanie) jest zakodowane w pojedynczych słowach stanowiących dany tekst.

- **Podejście bazujące na wzorcach (*pattern-based approach*)** – w tym podejściu istotne jest przekonanie, że nacechowanie opinii wyznaczają nie pojedyncze słowa, ale zbudowane z nich frazy/związki frazeologiczne. Tak więc konieczne jest wyszukanie wśród słów związków wyrazowych.
- **Podejście bazujące na ontologiach (*ontology-based approach*)** – pojedyncza opinia dotycząca produktu lub usługi może zostać przedstawiona jako instancja ontologii. Następnie instancje te mogą zostać porównane, a na tej podstawie reprezentowane przez nie opinie mogą zostać zaklasyfikowane do jednej z utworzonych grup.
- **Podejście oparte na uczeniu maszynowym (*machine learning approach*)** – dzięki zastosowaniu uczenia maszynowego można zbudować system, który nie tylko na podstawie odpowiednio dobranego uczącego zbioru opinii będzie je klasyfikował do odpowiednich grup, ale również będzie się rozwijał wraz z pojawieniem się nowych, specyficznych opinii.



**Rys. 1.** Wykorzystanie różnych podejść text miningowych do poszczególnych działań w ramach automatycznej analizy opinii konsumenckich

Źródło: opracowanie własne.

W modelu opartym na uczeniu maszynowym stosowane klasyfikatory wykorzystują odległość pomiędzy opiniami. Jest ona podstawą przypisania opinii pozytywnego lub negatywnego nacechowania.

Jednym ze sposobów klasyfikacji metod pomiaru podobieństwa dokumentów tekstowych jest klasyfikacja ze względu na fakt wykorzystania wiedzy dziedzinowej. Wyróżnić tu można metody niekorzystające z wiedzy dziedzinowej (głównie oparte na macierzy częstości) i te wykorzystujące wiedzę dziedzinową (głównie w postaci ontologii). W niniejszej pracy wykorzystana zostanie miara oparta na macierzy częstości.

### 3. Materiały i metody badań

Opinie to specyficzny rodzaj danych tekstowych, które mają subiektywny charakter – wyrażają stosunek autora wypowiedzi do przedmiotu opinii. W niektórych serwisach opinie słowne są wspierane oceną punktową lub gwiazdkami. Opinie można podzielić na grupy według ich formatu [Liu 2007]:

- format 1: zalety i wady oraz podsumowanie,
- format 2: zalety i wady,
- format 3: dowolny.

Analiza doniesień literaturowych [Grabowski 2012; Abramowicz 2008] pozwala wyodrębnić dwa główne podejścia stosowane przy ocenie miar podobieństwa dokumentów tekstowych (w tym opinii):

a) podejście polegające na określeniu zbieżności pomiędzy podobieństwem (lub odległością) wyznaczonym za pomocą ocenianej miary a podobieństwem rzeczywistym między dokumentami (określonym na przykład przez człowieka lub wynikającym ze sposobu generowania dokumentów na potrzeby badań symulacyjnych),

b) podejście polegające na analizie własności stosowanych miar poprzez analizę teoretyczną lub symulacyjną.

Przy ocenie zbieżności pomiędzy podobieństwem (lub odległością) wyznaczonym za pomocą ocenianej miary a podobieństwem rzeczywistym pomiędzy dokumentami zauważyć można jedną bardzo istotną trudność. Jest nią dostępność oceny rzeczywistego podobieństwa pomiędzy badanymi tekstami.

Przy dużych zbiorach opinii nieefektywne jest wyznaczanie rzeczywistego podobieństwa badanych tekstów przez badacza lub eksperta w danej dziedzinie. W badaniu przyjęto zatem założenie, że dwie opinie są tym bardziej do siebie podobne, im bardziej zbieżna jest ich ocena punktowa. Stanowi to punkt odniesienia do oceny przydatności miar podobieństwa dokumentów tekstowych, wyznaczając rzeczywiste podobieństwo pomiędzy badanymi tekstami.

Badania empiryczne podzielone zostały na sześć etapów:

1. ekstrakcja opinii,
2. utworzenie macierzy częstości,
3. wygenerowanie macierzy odległości opartej na ocenach użytkowników,
4. wygenerowanie macierzy odległości opartej na macierzy częstości,
5. porównanie dwóch macierzy odległości,
6. użycie algorytmu genetycznego.

W badaniach wykorzystano pakiety tm, RODBC, proxy i genalg języka R oraz aplikację RapidMiner.

### 3.1. Ekstrakcja opinii i macierz częstości

Opinie wykorzystane w badaniu pochodziły z serwisu Booking.com, pozwalającego na rezerwację pokoi hotelowych i udostępniającego opinie gości o hotelach. Wybrane zostały jedynie opinie w języku polskim. Podzielono je według hoteli, których dotyczyły. Wszystkie opinie w serwisie Booking.com dostępne są w formacie 2 (wady i zalety hotelu). Ponadto do każdej opinii dołączona jest punktowa ocena hotelu.

Pierwszym krokiem w badaniu była ekstrakcja danych z serwisu do bazy danych. W tym celu wykorzystano własny skrypt napisany w języku PHP oraz baza MySQL. Następnie dla każdej opinii polskojęzycznej w bazie danych wygenerowano plik tekstowy zawierający jedynie wady i zalety ocenianego hotelu. W nazwie pliku pojawiała się punktowa ocena. Pliki zostały zgrupowane w folderach zawierających nazwę ocenianego hotelu.

W drugim kroku na podstawie wygenerowanych plików tekstowych utworzono macierz częstości (macierz, której kolumny reprezentują dokumenty, a wiersze – wyrazy). Dla opinii o każdym hotelu tworzone były osobne macierze częstości.

### 3.2. Macierze podobieństwa i ich porównanie

Kolejnym krokiem w badaniu empirycznym było utworzenie dwóch macierzy odległości. Jedna z nich opierała się na ocenach punktowych. Dla każdej pary opinii policzono różnicę wartości ocen punktowych. Pracę tę wykonano w arkuszu kalkulacyjnym, a następnie w postaci gotowej macierzy odległości wczytano do programu R przy wykorzystaniu pakietu RODBC.

W dalszej części badań na podstawie macierzy częstości utworzono macierz odległości z wykorzystaniem odległości kątovej [Deza i Deza 2009]. W tym celu wykorzystany został pakiet tm języka R. Miarę tę wybrano ze względu na dobre wyniki osiągane przy jej wykorzystaniu we wcześniejszych analizach [Wójcik 2011].

	row.names	opinion1(8.8).txt	opinion10(4.6).txt	opinion11(7.9).txt	opinion12(8.8).txt	opinion13(7.1).txt
1	opinion1(8.8).txt	3.986534	3.511385	4.280387	3.986534	4.291590
2	opinion10(4.6).txt	4.569674	3.917934	4.841894	4.569674	4.829189
3	opinion11(7.9).txt	3.459496	2.992120	3.725890	3.459496	3.735925
4	opinion12(8.8).txt	4.188922	3.702500	4.491796	4.188922	4.503185
5	opinion13(7.1).txt	4.094811	3.560351	4.360887	4.094811	4.341064
6	opinion14(6.7).txt	3.379627	2.880257	3.639425	3.379627	3.610230
7	opinion15(6.7).txt	4.135343	3.590580	4.403882	4.135343	4.391938
8	opinion16(5.4).txt	4.186713	3.571949	4.447880	4.186713	4.427411
9	opinion17(7.5).txt	4.388495	3.855891	4.656962	4.388495	4.650412
10	opinion18(7.9).txt	4.126182	3.607475	4.391742	4.126182	4.402865

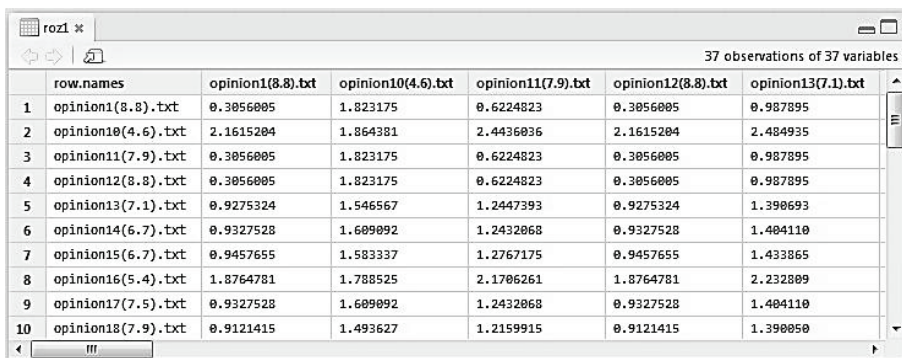
Rys. 2. Fragment macierzy odległości pomiędzy macierzami odległości opartymi kolejno na ocenach użytkowników i na macierzy częstości dla pojedynczego hotelu

Źródło: opracowanie własne w RStudio.

Po utworzeniu obydwu macierzy odległości (opartej na ocenach użytkowników i na macierzy częstości) zostały one porównane (policzona została odległość pomiędzy nimi). Wykorzystano do tego gotową funkcję *dist* w języku R z pakietu *proxy*. Rysunek 2 przedstawia fragment wyników tego porównania.

### 3.3. Algorytm genetyczny

W kolejnym etapie badań wykorzystano algorytm genetyczny [Goldberg 1989], którego zadaniem jest dobór zestawu wag używanych do ważenia macierzy częstości (wagi mogą przyjmować wartości 0 lub 1) tak, aby macierz odległości oparta na niej była jak najbardziej podobna do macierzy odległości opartej na ocenach użytkowników. Minimalizowana jest macierz będąca wynikiem działania funkcji *dist* dla obu macierzy odległości. Funkcja *dist* służy do wyznaczania podobieństwa pomiędzy macierzami.



The screenshot shows an RStudio window with a data frame containing 10 rows and 7 columns. The columns represent different opinion files, and the rows represent the pairwise distances between them. The values are numerical, ranging from approximately 0.9 to 2.4.

row.names	opinion1(8.8).txt	opinion10(4.6).txt	opinion11(7.9).txt	opinion12(8.8).txt	opinion13(7.1).txt
1 opinion1(8.8).txt	0.3056005	1.823175	0.6224823	0.3056005	0.987895
2 opinion10(4.6).txt	2.1615204	1.864381	2.4436036	2.1615204	2.404935
3 opinion11(7.9).txt	0.3056005	1.823175	0.6224823	0.3056005	0.987895
4 opinion12(8.8).txt	0.3056005	1.823175	0.6224823	0.3056005	0.987895
5 opinion13(7.1).txt	0.9275324	1.546567	1.2447393	0.9275324	1.390693
6 opinion14(6.7).txt	0.9327528	1.609092	1.2432068	0.9327528	1.404110
7 opinion15(6.7).txt	0.9457655	1.583337	1.2767175	0.9457655	1.433865
8 opinion16(5.4).txt	1.8764781	1.788525	2.1706261	1.8764781	2.232809
9 opinion17(7.5).txt	0.9327528	1.609092	1.2432068	0.9327528	1.404110
10 opinion18(7.9).txt	0.9121415	1.493627	1.2159915	0.9121415	1.390050

**Rys. 3.** Fragment macierzy odległości pomiędzy macierzami odległości opartymi kolejno na ocenach użytkowników i na macierzy częstości dla pojedynczego hotelu, powstałej w wyniku działania algorytmu genetycznego

Źródło: opracowanie własne w RStudio.

Wykorzystany w tym celu został wektor wag o liczbie elementów równej liczbie słów w macierzy częstości. Algorytm w 100 powtórzeniach dokonywał zmian w wektorze wag, następnie ważono macierz częstości i porównywano macierze odległości. Rysunek 3 przedstawia fragment macierzy z ostatecznym wynikiem dla jednej z podjętych prób (w każdej z prób uzyskiwano nieznacznie różniące się wyniki).

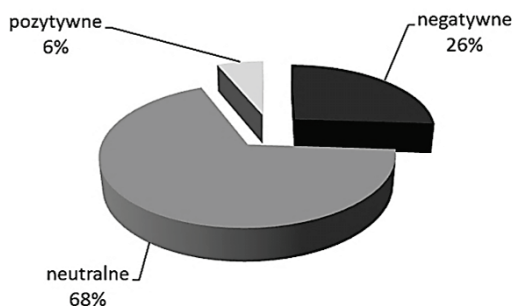
## 4. Wyniki badań empirycznych

Eksperyment przeprowadzony został na opiniach dotyczących pięciu hoteli. Zbiory opinii liczyły odpowiednio 18, 32, 37, 47 oraz 53 opinie w formie dokumentów tekstowych. Były to wszystkie opinie w języku polskim na temat danych hoteli

dostępne w serwisie Booking.com. Kryterium wyboru hoteli do badań była różnorodność opinii internautów na ich temat oraz ogólna ocena punktowa hoteli. Na jednym ze zbiorów danych badanie zostało wykonane pięciokrotnie. Był to zbiór liczący 37 opinii. Został on wybrany jako najbardziej reprezentatywny.

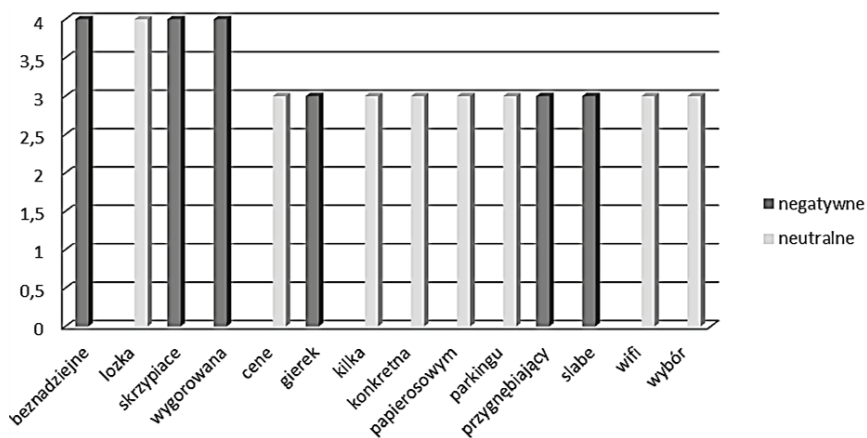
Celem analizy było przetestowanie narzędzia, jakim jest algorytm genetyczny, i ocena jego przydatności w automatycznej analizie opinii konsumenckich. Z punktu widzenia automatycznej analizy opinii konsumenckich ważną kwestią jest nacechowanie słów istotnych.

Zaczynając od zagregowanych wyników z pięciu prób przeprowadzonych dla pojedynczego hotelu na rysunku 4 można zauważyć, że 68% słów istotnych miało nacechowanie neutralne. Pomędzy pozostałymi słowami widoczna jest dysproporcja na korzyść słów negatywnych (26% w porównaniu z 6% słów pozytywnych). Odzwierciedla to ogólną ocenę danego hotelu przez jego gości.



**Rys. 4.** Nacechowanie słów istotnych, wybranych w pięciu powtórzeniach badania na zbiorze opinii o pojedynczym hotelu

Źródło: opracowanie własne.



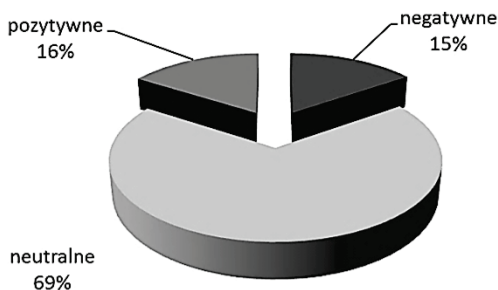
**Rys. 5.** Słowa powtarzające się w pięciu powtórzeniach badania na zbiorze opinii o pojedynczym hotelu

Źródło: opracowanie własne.

Z kolei rysunek 5 przedstawia słowa istotne, które powtarzały się w kolejnych eksperymentach. Żadne ze słów nie wystąpiło we wszystkich pięciu próbach. Pojawiały się jednak takie, które wystąpiły cztero- lub trzykrotnie. Warto zauważyć, że są to słowa neutralne i negatywne. Zaledwie dwa słowa pozytywne powtórzyły się dwukrotnie (nie zostało to ujęte na wykresie).

Na uwagę zasługuje fakt, iż słowa istotne z punktu widzenia celu analizy nie są słowami występującymi najczęściej w opiniach na temat danego hotelu.

Wykres na rys. 6 przedstawia wyniki pięciu eksperymentów przeprowadzonych na różnych zbiorach danych. Ponownie słowa neutralne stanowią większość wśród słów istotnych (69%). Wartość ta jest podobna do uzyskanej w przypadku badania na jednym zbiorze danych. W tym przypadku można jednak zaobserwować równowagę pomiędzy słowami pozytywnymi i negatywnymi (odpowiednio 16% i 15%) z nieznaczną przewagą tych pierwszych.



**Rys. 6.** Nacechowanie słów istotnych wybranych w pięciu powtórzeniach badania na zbiorach opinii o różnych hotelach.

Źródło: opracowanie własne.

Dominacja słów o nacechowaniu neutralnym, głównie rzeczowników, wśród słów istotnych sugeruje konieczność powiązania ich z wyrazami pozytywnymi lub negatywnymi, obrazującymi nastawienie autora wypowiedzi do poszczególnych aspektów przedmiotu opinii. Możliwość taką daje wykorzystanie podejścia bazującego na wzorcach, które pozwala na identyfikację w tekście związków frazeologicznych.

## 5. Podsumowanie

W artykule przedstawione zostały wyniki badań symulacyjnych, dotyczących doboru optymalnego zestawu słów istotnych podczas obliczania podobieństwa (lub odległości) dokumentów tekstowych, stosowanego na potrzeby automatycznej analizy opinii konsumenckich. Zastosowanie algorytmu genetycznego pozwoliło na efektywną selekcję zestawów słów istotnych.

Wśród słów istotnych dominują te o nacechowaniu neutralnym. O przewadze słów pozytywnych nad negatywnymi (lub na odwrót) decyduje sumaryczne nacechowanie opinii na temat danego dobra lub usługi. Ograniczenie liczby słów do tych istotnych nie wpływa na ocenę nacechowania opinii, lecz przyczynia się do zwiększenia efektywności badań poprzez ograniczenie liczby analizowanych danych.

Przeprowadzone analizy pozwalają stwierdzić, iż prowadzenie dalszych badań jest uzasadnione. W przyszłości są więc planowane:

- zwiększenie zbioru badanych opinii,
- zastosowanie podejścia opartego na wzorcach do identyfikacji związków wyrazowych występujących w opiniach,
- próba identyfikacji zestawów słów istotnych dla ogółu opinii na dany temat (np. dla wszystkich opinii o hotelach),
- próby wykorzystania wyników badań w dalszych analizach.

## Literatura

- Abramowicz W. (2008), *Filtrowanie informacji*, Wydawnictwo Akademii Ekonomicznej w Poznaniu, Poznań.
- Deza M.M., Deza E. (2009), *Encyclopedia of distances*, Springer-Verlag, Berlin – Heidelberg.
- Goldberg D.E. (1989), *Genetic Algorithms in Search, Optimization, and Machine Learning*. Reading, Addison-Wesley, Massachusetts.
- Grabowski M. (2012), *Naukowa legitymizacja obszaru Systemów Informacyjnych Zarządzania*. Wydawnictwo Uniwersytetu Ekonomicznego w Krakowie, Kraków.
- Liu B. (2007), *Web DataMining. Exploring Hyperlinks, Contents, and Usage Data*, Springer-Verlag, Berlin – Heidelberg.
- Lula P., *Automatyczna analiza opinii konsumenckich*, [w:] *Klasyfikacja i analiza danych – teoria i zastosowania*, red. K. Jajuga, M. Walesiak, Prace Naukowe Uniwersytetu Ekonomicznego we Wrocławiu nr 176, Taksonomia 18, Wydawnictwo Uniwersytetu Ekonomicznego, Wrocław, s. 53-62.
- Lula P., Wójcik K. (2011), *Sentiment analysis of consumer opinions written in Polish*, „Economics and Management” no 16, s. 1286-1291.
- Ohana B., Tierney B. (2009), *Sentiment Classification of Reviews Using SentiWordNet*, IT&T Conference, Dublin Institute of Technology, Dublin.
- Pang B., Lee L. (2008), *Opinion Mining and Sentiment Analysis*, „Foundations and Trends in Information Retrieval” 2(1-2), s. 1-135.
- Wójcik K. (2011), *Analiza porównawcza miar podobieństwa tekstów*, [w:] *Klasyfikacja i analiza danych – teoria i zastosowania*, red. K. Jajuga, M. Walesiak, Prace Naukowe Uniwersytetu Ekonomicznego we Wrocławiu nr 176, Taksonomia 18, Wydawnictwo Uniwersytetu Ekonomicznego, Wrocław, s. 340-348.
- Wikipedia. <http://pl.wikipedia.org/wiki> (12 października 2012).



## SELECTION OF THE OPTIMAL SET OF RELEVANT WORDS IN CONSUMERS OPINIONS IN THE CONTEXT OF THE OPINION MINING

**Summary:** Sentiment analysis is the research area that can have a significant impact on today's business. A lot of sentiment analysis studies base on similarity (or distance) between opinions as special kinds of texts. The problem that must be solved is to determine the best similarity measure to deal with opinions. The main objective of this paper is to conduct the analysis of the selection of optimal set of words relevant when calculating the similarity (distance) of text documents in the context of the opinion mining. Computations are conducted in R language, RapidMiner application and spreadsheet.

**Keywords:** text-mining, Web-mining, taxonomy, text document classification, evaluation of measures.