**Paper**

# A *k*-Nearest Neighbors Method for Classifying User Sessions in E-Commerce Scenario

Grażyna Suchacka[1], Magdalena Skolimowska-Kulig[1], and Aneta Potempa[2]

[1] Institute of Mathematics and Informatics, Opole University, Opole, Poland
[2] GEA Technika Cieplna Sp. z o.o., Opole, Poland

**Abstract—This paper addresses the problem of classification of user sessions in an online store into two classes: buying sessions (during which a purchase confirmation occurs) and browsing sessions. As interactions connected with a purchase confirmation are typically completed at the end of user sessions, some information describing active sessions may be observed and used to assess the probability of making a purchase. The authors formulate the problem of predicting buying sessions in a Web store as a supervised classification problem where there are two target classes, connected with the fact of finalizing a purchase transaction in session or not, and a feature vector containing some variables describing user sessions. The presented approach uses the *k*-Nearest Neighbors (*k*-NN) classification. Based on historical data obtained from online bookstore log files a *k*-NN classifier was built and its efficiency was verified for different neighborhood sizes. A 11-NN classifier was the most effective both in terms of buying session predictions and overall predictions, achieving sensitivity of 87.5% and accuracy of 99.85%.**

*Keywords—data mining, e-commerce, k-Nearest Neighbors, k-NN, log file analysis, online store, R-project, supervised classification, Web mining, Web store, Web traffic, Web usage mining.*

## 1. Introduction

Electronic commerce has gained tremendous popularity in recent years, especially in the area of Business-to-Consumer (B2C) trade, completed through Web stores. A still increasing number of online customers and their heterogeneous needs motivate online retailers to search for ways of predicting e-customer behavior and providing them with customized service. Especially valuable would be the ability to identify potential or future buyers as it could make it possible to enforce different personalized service strategies, e.g., by offering special discounts to encourage undecided visitors to buy or by presenting them with complementary products to increase online sales.

The problem of analyzing Web user behavior and anticipating their needs has been intensively explored. Various approaches have been proposed in the literature to discover hidden patterns in user navigation paths and online purchase transactions [1]–[4], as well as to identify users with high purchasing intentions and predict online sales [5]–[9]. This study fits into this research area by addressing the

problem of predicting online purchases. Authors propose an approach to classify user sessions in online store into buying sessions and browsing sessions based on some session features using a *k*-Nearest Neighbor (*k*-NN) technique. The rest of the paper is organized as follows. Section 2 provides a brief background on *k*-NN classification. Section 3 discusses related work on application of *k*-NN classifiers in e-commerce environment. Section 4 presents a methodology underlying the presented approach and Section 5 discusses the performance evaluation of the approach using data from a real online store. The final Section 6 provides some concluding remarks and directions for future work.

## 2. *K*-Nearest Neighbors Classifier

A *k*-Nearest Neighbors algorithm is a supervised learning technique often used in pattern recognition for classification, although it can also be used for estimation and prediction [10]. A *k*-NN classifier is memory-based (instance-based) and requires no model to be fitted, it is also conceptually simple. No explicit training procedure is required for the set of observations apart from collecting vectors of features with labels of classes they belong to. All intensive computations are performed at classification, which involves two steps for each test observation: finding $k$ nearest neighbors among the observations in a training set and performing a majority voting among the retrieved $k$ neighbors to assign the most frequent class label among them [11].

Let $\{(x_i, y_i), \ i = 1, 2, \ldots, n\}$ be the training set of observations, where $x_i$ is a vector of features and $y_i$ is a class label of $x_i$. We assume that every $x_i$ is in some multidimensional feature space with a metric $\rho$ and $y_i \in \{1, 2, \ldots, l\}$, where $l$ is the number of considered classes, $i = 1, 2, \ldots, n$. The task is to assign an unlabeled vector $x$ to a proper target class from $\{1, 2, \ldots, l\}$. The simplest version of *k*-NN algorithm is the nearest neighbor rule (1-NN) which assigns $x$ to the class of its closest neighbor. It means that if $x_j$, for some $j \in \{1, 2, \ldots, n\}$, is the nearest to $x$ in the sense of distance $\rho$, i.e.:

$$x_j = \arg \min_{\{x_i, 1 \leq i \leq n\}} \rho(x, x_i), \qquad (1)$$

then the rule labels $x$ the number $y_i$. The *k*-NN method for $k > 1$ is a natural extension of the foregoing rule. It

classifies *x* by assigning it the label which is most frequently represented among the *k* nearest training points $x_i$, where *k* is a user-defined constant. Thus, a decision is made by examining class labels of the *k* nearest neighbors and taking a majority vote.

The 1-NN rule usually classifies with an error rate greater than the minimum possible one – the Bayes rate. However, when the number of observations *n* tends to infinity, the error rate is not worse than twice the Bayes rate. For more details and discussion on the *k*-NN classification see e.g. [12] and [13].

The most common similarity measure (distance measure) between observations is the Euclidean metric. The Euclidean distance between two *J*-dimensional vectors *a* and *b* is expressed by the formula:

$$d_{a,b} = \sqrt{\sum_{j=1}^{J} (a_j - b_j)^2}. \qquad (2)$$

Additional metrics which can be used as distance measures are, e.g., standardized Euclidean, weighted Euclidean, Mahalanobis, Minkovsky, and Chebyshev distances, or in the case of a discrete type of data – a Hamming distance.

Huge advantages of the *k*-NN algorithm, especially important in practical applications, are its great scalability, linear computational complexity, robustness to data sparseness and skewed target class distribution, and interpretable target class predictions [11]. The *k*-NN method has been successfully applied in a large number of classification problems, like handwriting detection, gene expression, EKG patterns or satellite image scene detection. It has also been applied to some problems related to WWW and e-commerce, which are briefly reviewed in the Section 3.

# 3. Related Work

The most popular application area of the *k*-NN method in the e-commerce environment has been recommender systems. Generally, product recommendation methods in online stores may be divided into content-based recommendation and collaborative filtering (CF) recommendation. Content-based recommendation methods are based on the similarity of products (item profiles are built). On the other hand, CF recommender systems try to assess the utility of items for a given customer based on the items connected with other, similar users (user profiles are built) [14].

A traditional CF method is similarity-based, i.e., products to be recommended to a particular user are selected based on product preferences of a group of their nearest neighbors (with preferences similar to those of the target user) [15]. A comprehensive survey of CF techniques, divided into memory-based, model-based, and hybrid CF algorithms have been discussed in [16].

A key task in CF recommender systems is computing the all-to-all similarity between customers in order to form the neighborhood for a particular customer, i.e., a group of the most similar customers in terms of preferences, tastes, and

purchase patterns. The proximity between two e-customers in conventional CF methods has been the most frequently measured using the Pearson correlation coefficient [15], [17]–[19] or the cosine metric [11], [17]. Many extensions to the standard correlation-based and cosine-based techniques have been proposed, including default voting, inverse user frequency, case amplification, and weighted-majority prediction [14].

The most popular method for neighborhood formation is a center-based scheme, that forms a neighborhood for a particular customer by simply selecting the *k* nearest other customers. Other methods may be applied as well, e.g., in the case of very sparse data sets an aggregate neighborhood scheme may be applied, which allows the nearest neighbors to affect the process of a neighborhood formation for a particular customer [17].

Two main techniques are generally used to determine the neighborhood size [18]. The first one is a *correlation-thresholding* technique, in which the nearest neighbors are customers with absolute correlates greater than a given threshold. The second technique is a *best-k-neighbors* technique, in which the best *k* correlates are selected for the neighbors. In practice, the value of *k* is often determined experimentally.

Many improvements have been proposed for CF recommendation in e-commerce. Some studies have used the fact that a low-dimensional space is less sparse than the corresponding high-dimensional space and thus, they have applied dimensionality reduction techniques, e.g., Singular Value Decomposition (SVD) to generate a low-dimensional feature space in which the neighborhood has then been formed [17]. In [18] a CF-based recommendation methodology based on Web usage mining and product taxonomy was proposed to address the sparseness and scalability problems of CF recommender systems. The product taxonomy was used to improve the performance of searching for nearest neighbors through dimensionality reduction of the rating database. Jiang *et al.* [19] proposed a clustering-based *k*-NN approach which combines the *k*-NN and the iterative clustering algorithm. The iterative clustering approach allows them to solve the data sparseness problem by fully exploiting the voting information first. Then, a cluster-based *k*-NN is applied to improve the performance of CF.

Other application areas of the *k*-NN method in the e-commerce environment besides recommender systems have been based on the Web contents analysis. In [20] a method for e-commerce website trust assessment based on the analysis of textual contents and page layout was proposed. Two approaches were applied to construct a feature vector for each analyzed document: in a *baseline* approach all words extracted from the document were used and in a *EC-word* approach the extracted words were mapped into the meaningful groups of e-commerce terminology. Then, three different methods were applied to classify the text into a proper trust level class: *k*-NN, Naïve Bayes, and Support Vector Machine (SVM) based on Sequential Minimal Op-

timization (SMO) training. The $k$-NN method was used to find $k$ nearest neighbors of the analyzed document in the training set of documents using the Euclidean distance as a similarity measure. The categories of the nearest neighbors were used to determine the resultant trust level class. Results of $k$-NN classification were good, although a bit worse than SVM results.

In [11] the $k$-NN method was used in a hierarchical approach to the problem of large-scale text-based item categorization on an e-commerce site (a major online auction site). Categorization of items in very large datasets was formulated as a supervised classification problem where the categories are the target classes and words composing some textual description of the items are the features. The classification problem was decomposed into a coarse level task and a fine level task. The coarse level classifier was responsible for classifying items into so-called latent groups. A simple and scalable $k$-NN classifier for the reduced feature space was applied at this level. The similarity between items was measured using the cosine metric and the item was assigned to the class with the majority votes of its $k$ nearest neighbors. On the other hand, the fine level classifier was responsible for assigning items to the right class (category) in a given latent group. To this end, an SVM classifier was applied for all latent groups for the original latent group feature space.

To the best of authors knowledge, the problem of $k$-NN classification of e-customer sessions in terms of a purchase confirmation has not been investigated so far. This work is a continuation of authors previous studies on the application of various data mining techniques (e.g., association rules [8] and SVM [9]) to predict online purchases.

# 4. Research Methodology

## 4.1. Analysis of the Website Contents and Distinguishing Session States

A typical online store is implemented as a website hosted on a Web server on the Internet. The website consists of many pages, each of which is related to some function (e.g., reading general information about the store, searching for a product, adding the product to a shopping cart, etc.). At any time multiple users may interact with the site by opening different pages and performing various functions, i.e., there may be many active user sessions on the server.

The presented research was based on data obtained from an online bookstore (the name of the store is not given in the paper due to a non-disclosure agreement). The analyzed website contained not only typical Web store pages but also related pages with multimedia entertainment contents, like movies, quizzes, games, etc.

The website contents was analyzed in detail and as a result each page was assigned to one of the following 15 session states: *Home* – the home page of the bookstore, *Information* – pages containing general information about the

store, *Entertainment* – pages with entertainment contents, *Shipping* – pages with information on shipping cost and terms, visited before starting the checkout process, *Shipping_checkout* – pages with shipping information visited during the checkout process, *Browse* – browsing interactions, *Search* – interactions connected with searching for products according to some keywords, *Details* – pages with product information, *Add* – adding a product to the shopping cart, *Register_success* – successful user registration, i.e., creating a new user account, *Register_try* – pages connected with the registration process, other than the successful user registration, *Login_success* – successful user logging into the site, *Log_off* – user logout, *Checkout_try* – pages connected with the checkout process, other than a purchase confirmation, *Checkout_success* – purchase confirmation operation, i.e., a successful finalization of a purchase transaction.

## 4.2. Source Data

A single user session in an online store is represented on the Web server as a sequence of HTTP requests sent to the server by a client (i.e., an Internet browser). All HTTP requests coming to the server are registered in a server access log. In the case of a popular NCSA combined log format the following data is written in a log file for each request: client IP address, identifier and username, for authentication, date and time stamp, HTTP method, version of HTTP protocol, Uniform Resource Identifier (URI) of the requested server resource, HTTP status code, the volume of data transferred for the HTTP request in response, a referrer that linked the user to the store website, and a user agent string with some information on the client browser.

In this research the data recorded from 1 to 30 April 2014 was used. All user session data set contained 39,000 observations. This set was divided into two subsets: a training set and a test set. The training set was created by drawing 26,000 observations (among which 146 observations contained a purchase confirmation operation). The remaining 13,000 observations (including 72 observations with a purchase confirmation operation) were in the test set.

## 4.3. Reconstruction and Description of User Sessions

Using a dedicated C++ program data was read from log files, merged, pre-processed and cleaned. Based on IP addresses and user agent strings of HTTP requests, request streams for individual users were distinguished. Assuming that intervals between each two subsequent requests of the same user in session do not exceed a 30-minute threshold [21]–[24] user sessions were then identified.

Each user session was described with 23 session features. The first group of features includes 15 elements connected with visits to the session states occurred in session. A variable $V_{Checkout\_success}$ is a boolean variable (equal to one if a purchase transaction was successfully finalized in session and zero otherwise). Other variables: $V_{Home}$, $V_{Information}$, $V_{Entertainment}$, $V_{Shipping}$, $V_{Shipping_{checkout}}$,

$V_{Browse}$, $V_{Search}$, $V_{Details}$, $V_{Add}$, $V_{Register\_success}$, $V_{Register\_try}$, $V_{Login\_success}$, $V_{Log\_off}$, $V_{Checkout\_try}$, are connected with numbers of visits to the corresponding session states occurred in session.

The second group of features includes the following 6 session characteristics:

- $V_{Requests}$, the number of HTTP requests in session,

- $V_{Transfer}$, the volume of data downloaded in session (in kilobytes),

- $V_{Pages}$, the number of pages visited in session,

- $V_{Duration}$, the session duration (in seconds),

- $V_{Time\_per\_page}$, the mean time per page (in seconds),

- $V_{Source}$, the source of the visit (a reference from natural or paid search engine results, an e-mail newsletter, or a social media site, internal reference from a page with entertainment content, or other source).

Two last session features, $V_{Is\_bot}$ and $V_{Is\_admin}$, are boolean variables indicating whether the session was performed by a Web bot or the website administrator, respectively.

### 4.4. Problem Formulation

The research goal is to predict user sessions that ended with purchases, so the sessions were classified depending on if a purchase transaction was finalized in session or not. That is why two session classes (browsing sessions and buying sessions) are considered, differing in the value of $V_{Checkout\_success}$ variable, which is a class label $y_i \in \{0, 1\}$, $i = 1, 2, \ldots, n$, in fact. A vector of features, $x_i$, contains 22 predictor variables corresponding to the remaining session features.

$K$-NN classifiers for different $k$ values were built based on the training set and their performance was evaluated based on the test set. This research was realized using a free software for statistical computing, R-project [25].

## 5. Results and Discussion

First, each of 22 predictor variables for all observations in the training set and the test set was standardized to have a mean equal to zero and a variance equal to one. After the standardization $k$-NN classifiers for $k = 1, 2, \ldots, 40$ were built based on the training set. In all cases the proximity between user sessions was measured by using the Euclidean metric.

Then, the classifiers performance was evaluated using the test set. A classification result for each observation may be a buying session, which is considered to be positive, or a browsing session, considered to be negative. Depending on the actual value of the variable $V_{Checkout\_success}$ (1 or 0) for each observation, the session classification may be true or false. Results of the classification may be expressed with the numbers of true and false positives and negatives

by comparing predicted classifications vs. actual classifications:

- true positives (TP) is the number of correctly classified buying sessions,

- true negatives (TN) is the number of correctly classified browsing sessions,

- false positives (FP) is the number of browsing sessions which were incorrectly classified as buying sessions,

- false negatives (FN) is the number of buying sessions incorrectly classified as browsing sessions.

As one can see in Figs. 1 and 2, the classification results differ depending on the number of nearest neighbors taken into account when determining a class label. The highest number of correctly classified buying sessions, i.e., true positives was 63 and this result was achieved for $k$ ranging from 11 to 15 (Fig. 1). Thus, the classifier with 11–15 nearest neighbors was the most effective in predicting online purchases. On the other hand, the highest number of correctly classified browsing sessions, i.e., true negatives was 1292 for $k$ equal to 5 (Fig. 2).
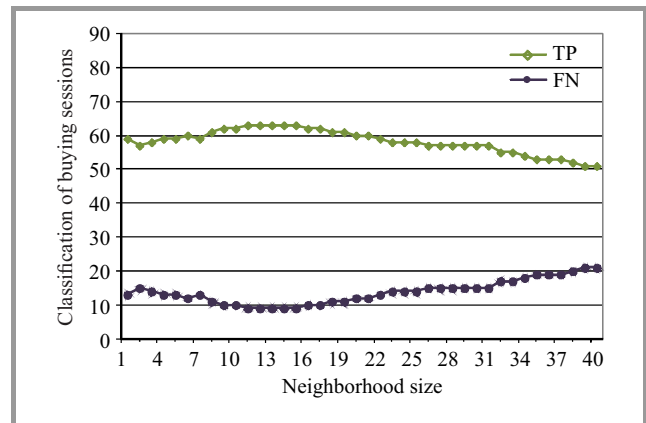


*Fig. 1.* Results of $k$-NN classification for buying sessions depending on the neighborhood size.
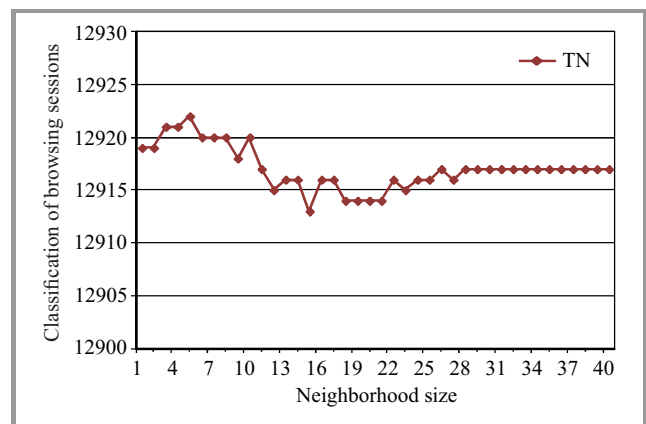


*Fig. 2.* Results of correct $k$-NN classification for browsing sessions depending on the neighborhood size.

The quality of the classifiers may be assessed with measures of predictive accuracy, error rate, and sensitivity. The accuracy is defined as the percentage of all correct classifications:

$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN}. \qquad (3)$$

Similarly, the error rate is the percentage of all incorrect classifications:

$$Error\ rate = \frac{FP+FN}{TP+TN+FP+FN}. \qquad (4)$$

A sensitivity measure is the percentage of correct classifications of buying sessions, so it is an estimate of the probability of predicting a buying session:

$$Sensitivity = \frac{TP}{TP+FN}. \qquad (5)$$

This measure is of the most importance from a retailer's point of view, because the ability to identify potential buyers in an online store makes it possible to apply different incentives to purchase, e.g., by offering special discounts to make an undecided visitor buy.

The performance rates for $k$-NN classifiers are illustrated in Fig. 3. All the classifiers have a very high predictive accuracy exceeding 99.74% and a correspondingly low error rate below 0.25%. The sensitivity of the classifiers is more differentiated, however, and ranges from 70.83% for $k$ equal to 40 to 87.50% for $k$ ranging from 11 to 15. One can observe that in general, as $k$ increases the classifier sensitivity increases as well until the neighborhood size of 15 and then continues to drop.
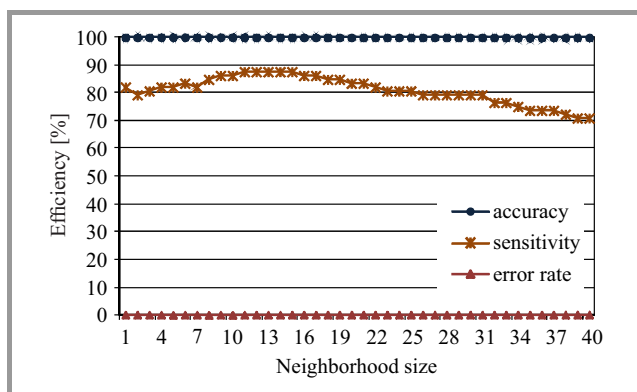


**Fig. 3.** The efficiency of $k$-NN classification in terms of all user sessions and buying sessions depending on the neighborhood size.

To sum up, the classifier taking into account the 11 nearest neighbors was the most effective in predicting buying sessions and in terms of overall predictions, achieving sensitivity of 87.5% and accuracy of 99.85%.

## 6. Conclusions

In this paper the problem of supervised classification of user sessions in a Web store in terms of purchase transaction realization was investigated. Using data from a real

online retailer, the authors described each user session with a 22-element feature vector and performed a $k$-NN classification of sessions into two classes: buying sessions and browsing sessions. Evaluation of the efficiency of $k$-NN classifiers for different neighborhood sizes showed that a classifier based on the 11 nearest neighbors was the most effective, achieving the overall predictive accuracy of 99.85%, while being capable of correctly classifying 87.5% buying sessions.

For future work, the authors would like to continue research on $k$-NN classification of e-customer sessions by examining various similarity measures and voting schema among the $k$ nearest neighbors, as well as to apply various dimensionality reduction techniques to the session feature space. It would also be worth verifying the efficiency of the proposed approach for other e-commerce data sets.

## References

[1] Q. Duan, J. Li, and Y. Wang, "The application of fuzzy association rule mining in e-commerce information system mining", *Adv. Engin. Forum*, vol. 6–7, pp. 631–635, 2012.

[2] G. Kuang and Y. Li, "Using fuzzy association rules to design e-commerce personalized recommendation system", *TELKOMNIKA Indonesian J. Elec. Engin.*, vol. 12, no. 2, pp. 1519–1527, 2014.

[3] Y.-S. Lee and S.-J. Yen, "Mining web transaction patterns in an electronic commerce environment", in *Advances in Web and Network Technologies, and Information Management – Proc. APWeb/WAIM'07 International Workshops*, Huang Shan, China, 2007, *LNCS*, vol. 4537, pp. 74–85. Springer, 2007.

[4] N. D. Thuan, N. G. Toan, and N. L. V. Tuan, "An approach mining cyclic association rules in e-commerce", in *Proc. 15th Int. Conf. Network-Based Inform. Syst. NBiS 2012*, Melbourne, Australia, 2012, pp. 408–411.

[5] W. Hop, "Web-shop order prediction using machine learning", Masters Thesis, Erasmus University Rotterdam, 2013.

[6] M. Mohammadnezhad and M. Mahdavi, "Providing a model for predicting tour sale in mobile e-tourism recommender systems", *IJITCS*, vol. 2, no. 1, pp. 1–8, 2012.

[7] N. Poggi *et al.*, "Web customer modeling for automated session prioritization on high traffic sites", in *Proc. User Modeling'07*, Corfu, Greece, 2007, LNCS, vol. 4511, pp. 450–454. Springer, 2007.

[8] G. Suchacka and G. Chodak, "Practical aspects of log file analysis for e-commerce", in *Proc. Computer Networks'13*, Lwówek Śląski-Brunów, Poland, 2013, CCIS, vol. 370, pp. 562–572. Springer, 2013.

[9] G. Suchacka, M. Skolimowska-Kulig, and A. Potempa, "Classification of e-customer sessions based on Support Vector Machine", in *Proc. 29th Eur. Conf. Model. Simul. ECMS'15*, Albena, Bulgaria, 2015, pp. 594–600.

[10] D. T. Larose, Discovering Knowledge in Data: An Introduction to Data Mining. Wiley, 2005.

[11] D. Shen, J.-D. Ruvini, and B. Sarwar, "Large-scale item categorization for e-commerce", in *Proc. 21st ACM Int. Conf. Inform. Knowl. Manag. CIKM'12*, Maui, HI, USA, 2012, pp. 595–604.

[12] R. O. Duda, P. E. Hart, and D. G. Stork, *Pattern Classification*. Wiley, 2000.

[13] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning. Data Mining, Inference, and Prediction*. Springer Series in Statistics, 2009.

[14] G. Adomavicius and A. Tuzhilin, "Toward the next generation of recommender systems: A survey of the state-of-the-art and possible extensions", *IEEE Trans. Knowl. Data Eng.*, vol. 17, no. 6, pp. 734–749, 2005.
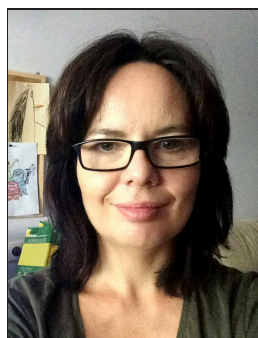
[15] J. Cho, K. Kwon, and Y. Park, "Collaborative filtering using dual information sources", *IEEE Intel. Syst.*, vol. 22, no. 3, pp. 30–38, 2007.

[16] X. Su and T. M. Khoshgoftaar, "A survey of collaborative filtering techniques", *Adv. Artif. Intel.*, vol. 2009, Article no. 4, 2009.

[17] B. Sarwar, G. Karypis, J. Konstan, and J. Riedl, "Analysis of recommendation algorithms for e-commerce", in *Proc. 2nd ACM Conf. Elec. Commerce EC'00*, Minneapolis, MN, USA, 2000, pp. 158–167.

[18] Y. H. Cho and J. K. Kim, "Application of web usage mining and product taxonomy to collaborative recommendations in e-commerce", *Expert Syst. Appl.*, vol. 26, no. 2, pp. 233–246, 2004.

[19] X.-M. Jiang, W.-G. Song, and W.-G. Feng, "Optimizing collaborative filtering by interpolating the individual and group behaviors", in *Proc. 8th Asia-Pacific Web Conf. Frontiers of WWW Res. Develop. APWeb'06*, Harbin, China, 2006, LNCS, vol. 3841, pp. 568–578. Springer, 2006.

[20] B. Soiraya, A. Mingkhwan, and C. Haruechaiyasak, "E-commerce web site trust assessment based on text analysis", *Int. J. Business Inform.*, vol. 3, no. 1, pp. 86–114, 2008.

[21] M. Adnan, M. Nagi, K. Kianmehr, R. Tahboub, M. Ridley, and J. Rokne, "Promoting where, when and what? An analysis of Web logs by integrating data mining and social network techniques to guide ecommerce business promotions", *Soc. Netw. Anal. Min.*, vol. 1, no. 3, pp. 173–185, 2011.

[22] L. D. Catledge and J. E. Pitkow, "Characterizing browsing strategies in the World-Wide Web", in *Proc. 3rd Int. World-Wide Web Conf. Technol., Tools Appl.*, 1995, pp. 1065–1073.

[23] Z. Chen, A. W.-C. Fu, and F. C.-H. Tong, "Optimal algorithms for finding user access sessions from very large Web logs", *World Wide Web*, vol. 6, no. 3, pp. 259–279, 2004.

[24] D. Stevanovic, N. Vlajic, and A. An, "Unsupervised clustering of Web sessions to detect malicious and non-malicious website users", *Procedia Comput. Sci.*, vol. 5, pp. 123–131, 2011.

[25] The R project for statistical computing [Online]. Available: http://www.r-project.org

Opole University, Poland. Her research interests include Web mining, Web analytics, and Quality of Web Service with special regard to electronic commerce.
E-mail: gsuchacka@math.uni.opole.pl
Institute of Mathematics and Informatics
Opole University
Oleska st 48
45-052 Opole, Poland



**Magdalena Skolimowska-Kulig** received her Ph.D. degree in Mathematics from University of Wrocław, Poland. Now she works as an assistant professor in the Institute of Mathematics and Informatics at Opole University, Poland. Her research interests include theory of probability and mathematical statistics.

E-mail: mskolimow@math.uni.opole.pl
Institute of Mathematics and Informatics
Opole University
Oleska st 48
45-052 Opole, Poland



**Aneta Potempa** received the B.S. degrees in Mathematics from Opole University, Poland and in Economy, accounting and audit from WSB Schools of Banking in Opole, Poland. She received the M.Sc. degree in Financial Mathematics from Opole University in 2014. She works in a controlling department in GEA Technika Cieplna, the company specializing in production of industrial heat exchangers.

E-mail: aneta.potempa@gea.com
GEA Technika Cieplna Sp. z o.o.
Kobaltowa st 2
45-641 Opole, Poland



**Grażyna Suchacka** received the M.Sc. degrees in Computer Science and in Management from Wrocław University of Technology, Poland. She received her Ph.D. degree in Computer Science from Wrocław University of Technology. Now she is an assistant professor in the Institute of Mathematics and Informatics at