

A Systematic Review on Web Page Classification

Ajose-Ismail, B.M¹, Osanyin Q.A²

¹Department of Computer Science, School of Applied Science,
The Federal Polytechnic Ilaro, Ogun State, Nigeria.

²Department of Computer Science, School of Applied Science,
The Federal Polytechnic Ilaro, Ogun State, Nigeria.

Abstract

With the increase in digital documents on the World Wide Web and an increase in the number of webpages and blogs which are common sources for providing users with news about current events, aggregating and categorizing information from these sources seems to be a daunting task as the volume of digital documents available online is growing exponentially. Although several benefits can accrue from the accurate classification of such documents into their respective categories such as providing tools that help people to find, filter and analyze digital information on the web amongst others. Accurate classification of these documents into their respective categories is dependent on the quality of training dataset which is dependent on the preprocessing techniques. Existing literature in this area of web page classification identified that better document representation techniques would reduce the training and testing time, improve the classification accuracy, precision and recall of classifier. In this paper, we give an overview of web page classification with an in-depth study of the web classification process, while at the same time creating awareness of the need for an adequate document representation technique as this helps capture the semantics of document and also contribute to reduce the problem of high dimensionality.

Keywords - Bags of words model, Classification, Machine learning, Document representation, TF-IDF, Web Page classification, LDA, Word2Vec.

I. INTRODUCTION

News and blogs webpages are today's most common sources for gathering information about current events. Information gotten from blogs and websites come in several categories in which users are only interested in certain topics within that category; for example business, entertainment, sports or politics. Aggregating and categorizing information from these sources seems to be a daunting task as the volume of digital documents available online is growing exponentially as a result of increased usage of the internet [1]. Automated web categorization is the key technology for this task. Web page classification

(WPC), also known as web page categorization, is the process of assigning a web page to one or more predefined category labels [2]. Web page classification problem can be divided into two categories: manual and automatic web page classification. Manual classification is a task that is performed by domain experts manually and it looks impractical because it will take lots of human effort and time [3]. While automatic web page classification is supervised machine learning problem where set of document is used to train the classifier, once training is done it is used to classify web pages [4]. While the former is tedious and time consuming, the latter saves lot of manpower and material resources and time [5]. Web classification is different from the standard text classification in some aspects: Traditional text classification is typically performed on structured documents which are stored in structured data stores such as relational databases and written with consistent styles which web collections do not possess [6], [7], [8]. Web pages are semi-structured documents formatted with HTML tags so that they may be rendered visually to users. Web documents also exist within a hypertext with connections within and outside the documents [6]. Several benefits can accrue from the accurate classification of documents into their respective categories such as providing tools that help people to find, filter and analyze digital information on the web. Also news filtering, document routing and personalization of information on the web are additional advantages that can be harvested from web page classification.

According to [9], the applications of WPC are as follows: Web directories provided by different search engines like Google, Yahoo amongst others can be constructed, maintained or expanded using advanced WPC techniques [10], [11]. WPC are used to improve the quality of search results. When a user types in a particular keyword, the numbers of relevant results are increased through WPC [12]. A question answer system uses WPC techniques to improve the quality of answers [13]. Web content filtering is another application of WPC [14]. Many WPC systems have been presented in literature over the years in which different perspectives have been taken to improve the performance of web classifiers

[15]. Machine Learning (ML) algorithms such as Naïve bayes, K-nearest neighbor, Decision tree, neural network, support vector machine and so on have been used previously by many researchers to achieve this task [16]. To achieve high classification result of the WPC system, an excellent representation of textual data (Preprocessing) should contain as much information as possible from the original document [17]. Also, the accuracy of most classification algorithms depends on the quality and size of training data which is dependent on the document representation technique [3]. The general problem of web classification can be divided in to three areas: document representation, classifier construction and classifier evaluation [18]. This paper provides an extensive study of web page classification process with a thorough review of feature selection techniques used in document representation using existing literature. The remainder of this paper is organized as follows: Section 2 proceeds with an overview of the web page classification process, highlights feature selection techniques used in document representation phase of the WPC. In Section 3, related works on web page classification. Finally the review is concluded in section 4, with also some future directions.

II. LITERATURE REVIEW

A. Web Page Classification Process

According to [16], the web page classification system is divided into several components as shown in Figure 1 below. The stages of the web page classification process includes: Creating a corpus of web pages, pre-processing / document representation, organization of the pre-processed pages, building the WPC model, obtaining a trained classifier, evaluating the classifier.

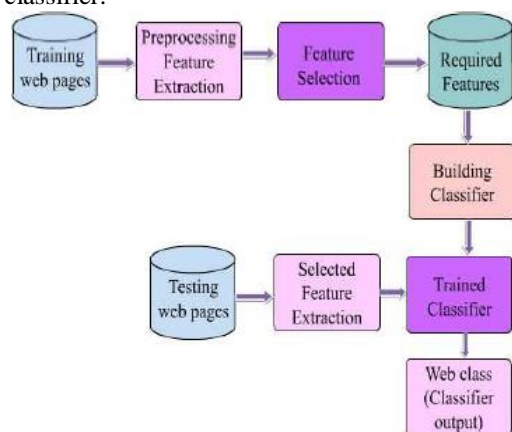


Figure 1: Web page Classification Process (Source: [19])

a) Corpus or Web Pages Training Dataset

The first stage in the web classification process proceeds with extracting the main contents of the webpage along with other web page elements such as

Internal and external hyperlinks, Metadata, Flash animation, Java script, Video Clips, Embedded objects, advertisement, Google ad-sense [2]. The extracted web contents are used in creating a corpus of labeled web pages i.e. training web pages which would be utilized by the classifier to building the learning system [6]. Already existing corpuses such as Reuters [3], [16], WebKb [4], [5], [9], Yahoo news dataset [20], sentiment Treebank [21], 20 News group dataset [17], imdb dataset [22] can be utilized for this process or by creating a custom made corpus (using a web crawlers) which can be used to automatically download web content [2].

b) Pre-processing / Document Representation

The next stage in the web page classification process is the pre-processing stage also known as Document Representation (DR) or dimensionality reduction in this context [11]. The pre-processing stage can be further divided into Feature Extraction (FE) and Feature Selection (FS) [18]. FE process begins by extracting the raw content of the pages and discard HTML tags and other WWW contents. Web page document are characterized by high dimensionality, the first technique to reduce this high dimensionality is FE [1], [4]. Then FE continues by performing tokenization (breaking a stream of text into words, phrases, symbols, or other meaningful elements called tokens), stemming / lemmatization [16]. After feature extraction, the next step in the pre-processing stage is Feature Selection (FS) which involves constructing a vector space model of the document to improve the scalability, efficiency and accuracy of a text classifier. The main idea of FS is to select a subset of features from the original documents [19]. Also with the inherent characteristics of web document which is high dimensional datasets, FS is used to reduces the feature space and improve the efficiency and accuracy of classifiers. Feature selection approaches can be broadly classified as filter, wrapper, and embedded. The most generic of all the approaches is the filter approach and it works irrespective of the data mining algorithm that is being used [3]. It typically employs measures like correlation, entropy, mutual information, and so forth which analyzes general characteristic of the data to select an optimal feature set. However, it is to be noted that wrapper and embedded methods often outperform filter in real data scenarios [19]. In the embedded approach, feature selection is a part of the objective function of the algorithm itself. Similar examples can be seen in decision tree, LASSO, LARS, 1-norm support vector, and so forth. In contrast to the above approaches which specifically select a subset of features, other techniques decomposes the original higher dimensional document-feature matrix into lower dimensional matrix, which effectively transforms the original feature space (determining the semantic relations between words, which defines the concept in the

document) [6]. Also, they are deficient in revealing inter-or intra-document statistical structure of the corpus [23]. Such methods include: bag of words model TF-IDF, Latent Semantic Indexing (LSI), Latent Semantic Analysis (LSA), Probabilistic Latent Semantic Indexing (PLSI), Word2Vec [17]. Each technique has its own pros and cons. Lots of discussions are ongoing in the pre-processing and document representation stage of the WPC system. DR is a very important step in web classification, because irrelevant and redundant features often degrade the performance of the classification algorithms both in speed and classification accuracy and also its tendency to reduce overfitting [19]. Also this stage has gained more attention recently than any other component of the WPC because effective dimension reduction makes the learning task more efficient and saves more storage space [24].

c) Obtaining the Required Features

The next stage after pre-processing is to gather the required feature set for classification which is usually achieved by creating matrix representation of the document vectors which would be fed to the classifier [19].

d) Building the WPC Model

After gathering the required features, the next stage is to build the WPC model using a classification algorithm with the selected features as the input data set. Several machine learning algorithm have been used for the building the model of the WPC system systems such as KNN, Support Vector Machine (SVM), Artificial Neural Network (ANN), Deep Learning [15] and so on. After training the classifier, the model obtained is then used to automatically classify new web pages to the appropriate category. Several authors have argued about the best ML technique for web page classification but literature has shown that the accuracy, generalization capabilities of any ML technique depends on the training data set i.e. the choice of the techniques used in the pre-processing stage have an overall effect on quality of the classifier [25].

e) Evaluating the Classifier

An evaluation measure is used to measure the performance of a WPC classifier. For each category T_n , a confusion matrix can be constructed as shown in the Figure 2 where ‘i’ denotes the number of true positive classifications, ‘j’ denotes the number of false positive classifications, ‘k’ denotes the number of false negative classifications and ‘l’ denotes the number of true negative classifications. For a perfect classifier j and k would both be zero [26].

TABLE 1: CONFUSION MATRIX (SOURCE: [26])

		Predicted	
		Class	
		T_n	Not T_n
Actual Class	T_n	I	j
	Not T_n	K	l

B. Feature Selection Techniques Used in WPC

According to Google index the volume of digital documents available online is over 130 trillion pages and its growing exponentially as a result of increased usage of the internet. Finding relevant and timely information from these documents are important for many applications [27]. The accuracy and generalization capabilities of the classifier in assigning a web page to its correct category is heavily dependent on the document representation [17]. Several authors have applied various DR techniques to improve the quality of the input dataset which inherently will increase the general performance of the WPC system. Each technique is fraught by one challenge or the other. Some of them are highlighted below:

a) Term Frequency-Inverse Document Frequency (TF-IDF)

TF-IDF is a model for document representation that is often used in information retrieval. It is a model that evaluates how important a word is to a document. It weighs the important words increasingly based on how frequently they appear in the document but decreases the weight proportionally as it occurs in other documents. TF-IDF can represent a document well by removing stop words from the documents. Some of the drawbacks of tf-idf is that it does not capture semantic similarity, does not respect word order and it is an unordered collection of words (Turney & Pantel, 2010).

b) Latent Semantic Indexing (LSI)

LSI is a popular information retrieval method that uses linear algebraic indexing method to produce low dimensional representations by word co-occurrence [28]. It utilizes the Singular Value Decomposition (SVD) algorithm on the sparse TF-IDF vector matrix, to create a denser matrix that approximately models the original document. It composes frequencies of terms as a term-document matrix. The term document matrix it is a sparse matrix whose rows correspond to terms and whose columns correspond to documents. LSI was used to solve the synonym and polysemy problem of TF-IDF. However, a major drawback of LSI are that, it does not capture multiple meanings of a word and it does not respect word order Also, LSA assumes that documents and features form a joint Gaussian model, while a Poisson distribution is typically observed and the resulting dimensions

might be difficult to interpret (Zhang, Yoshida & Tang, 2011).

c) Probabilistic Latent Semantic Indexing (PLSI)

To overcome some of the afore-mentioned problems with LSI, [29] introduced Probabilistic LSA (PLSA), which is a generative, graphical model enhancing latent semantic indexing (LSI) by a sounder probabilistic model. PLSI models the probability of each co-occurrence as a mixture of conditionally independent multinomial distributions. It uses EM Algorithm for its learning. PLSI is usually viewed as a more sound method as it provides a probabilistic interpretation, whereas LSI achieves the factorization by using only mathematical foundations. The core of PLSA is a statistical model which has been called aspect model. Although PLSI had promising results, it suffers from two limitations: the number of parameters is linear in the number of documents, and it is not possible to make inference for unseen data [23].

d) N-GRAM MODEL

The N-gram model is one of the most widely used models for feature representation. It assumes that the probability of a given word is only conditional on its preceding n-1 word, where n could be 1 (the unigram model), 2 (the bigram model), 3 (the trigram model), or any whole number. This approach converts a collection of text documents into feature vectors by recording the n-gram frequency counts, and uses the vectors as input to classifiers. According to (Elberrichi & Aljohar, 2007), some of the major strengths of N-GRAM are: no need to perform word segmentation, automatic capture of the roots of the most frequent words, independence towards the document language, tolerances with the spelling mistakes and the deformations. In addition, no dictionary or language specific techniques are needed. Also, N-GRAM suffers from data sparsity and high dimensionality [30].

e) Latent Dirichlet Analysis (LDA)

LDA is a probabilistic topic model that generates topics based on word occurrences from a corpus or set of documents [31]. It assumes documents are a blend of several topics and that each word in the document can be grouped under the document's topics. LDA is particularly useful for finding reasonably accurate mixtures of topics within a given document set. LDA is an unsupervised language model that transforms words from bag of words counts into continuous representative matrix. Also, LDA uses an unsupervised learning function which depends on words in the corpus which will determine the matching degree and thus will suffer from vocabulary mismatching problem [32].

f) Word2Vec

Word2Vec is a Recurrent Neural Network based implementation that can learn word embedding's. The 2 main architectures are CBOW (Continuous Bag-of-word) and Skip-gram (Continuous Skip-gram Model). CBOW tries to predict words from the context of words while skip-gram tries to predict the context from the words. In the CBOW model each input vector $w(t)$ is a column in the Matrix W . The CBOW model predicts a word $w(t)$ utilizing the context $w(t-n) \dots w(t-1), w(t+1) \dots, w(t+n)$, while the Skip-gram model predicts each word in the context utilizing the word $w(t)$. The Word2Vec framework aims at predicting the context of word or word based on their context. The word embedding's are learned through maximizing the objective function. With these word embedding's it can capture distributed representations of text to capture similarities among concepts [33] which is one of the major advantages of Word2Vec. However, a major drawback of word2Vec is that it does not model the global relationship between document to topics (Wang, Ma & Zhang, 2016).

III. RELATED WORKS ON WEB PAGE CLASSIFICATION

In the works of [1], they proposed a method to accurately and automatically classify web pages into different categories viz three phases: feature extraction, information learning and classification. In the methodology adopted, term document matrix is created using tf-idf, then the terms are used to extract object based features. Decision tree algorithm is then used to extract rules from the features extracted. The web pages are then classified using optimal firefly algorithm based Naive Bayes Classifier (FA-NBC) using the rules extracted. The proposed method was applied to WebKB datasets. Experimental results shows that their proposed method outperforms earlier methods such as KNN. Drawbacks identified in their work include: using tf-idf to construct the term document matrix does not capture any semantic similarity or form of grammatical analysis [17].

[34] proposed a method to analyze and categorize e-commerce websites automatically. In their methodology, e-commerce website were crawled, text preprocessing and the terms of the document were derived using tf-idf. The proposed method was applied to 1312 e-commerce and 1077 non e-commerce web site, preprocessing of the webpages, term weighted with tf-idf and classified using SVM. Experimental results shows that the produced method outperforms pure TF-IDF. Also the results shows a substantial increase in the accuracy of the classifier. A major gap identified in their work is that bag of words model like TF-IDF does not capture semantic similarity and respect word order of the document being represented [25].

In the works of [17], they proposed the use of a hybrid strategy that consist of Latent Dirichlet Allocation (LDA) and Word2Vec for document representation. Word2Vec create a vector representation of the document which shows the semantic relationship between the words of the document. Euclidean distance was used to measure and interpret similarity between document and topic in sparse space. Their methods was applied to 20 News group data using SVM classifier. Results obtained shows that their proposed methods outperforms earlier methods such as TF-IDF+ SVM, Word2Vec + SVM, LDA + SVM. One of the major drawback of their method is that improper calibration of the LDA parameters (e.g. number of topics, hyper-parameters), could potentially lead to sub-optimal results as most of the parameters for the LDA are imported from natural language community.

[15] proposed the categorization of web News using word2vec and deep learning. Earlier methods of automatic classification used supervised and unsupervised algorithm such as SVM, Naive Bayes and clustering respectively but are marred with several issues. The former can only handle supervised data but actual texts in web are not supervised data. The latter defines category automatically. Also, another problem with web classification is transforming the text in to constant dimensions. In their work, Word2Vec was used to train the news corpus into vectors. After obtaining the word vectors, pre-training of the vectors using autoencoder and training of the dataset using deep learning framework. Experimental test were carried out on web news site (Yahoo) containing 1,728,942 records and results obtained shows that deep learning produced an a better result than Naive Bayes but perform badly on training time.

IV. OBSERVATIONS

Representation of the input data (DR) is a crucial issue in web page classification and text classification systems at large. Several feature selection techniques have been proposed to solve the issue of semantic matching of unstructured data, but are marred with one issue or the other. Recently, there has been an increase in the use of SVM and KNN for text classification [26]. Also from extant literature, SVM, KNN and Naïve bayes are one of the most widely used ML algorithm for text classification [1], [5], [24]. In the work of [3], they decided to investigate this issue and compared SVM, KNN and Naïve Bayes on text classification tasks. Results obtained shows that SVM was not a clear winner, despite quite good overall performance. If a suitable pre-processing is applied to KNN and Naïve Bayes theory, these algorithms will achieve very good results and scales up to the performance of SVM. In light of this, there is need for an adequate document representation technique to retrieve the semantics of a web document. Optimized document representation

techniques such as hybridizing neural network language models (Word2Vec) and topic model (LDA) or Word2Vec and TF-Idf with optimizing the parameters of LDA with search algorithms (such as GA) will provide better semantics of the document in WPC. This hybrid approaches has shown to perform better (obtain the semantic features) by harnessing the strength of the individual technique in the arrangement Word2Vec and LDA [17] or Word2Vec and TF-Idf [24]. Also, proper calibration of the parameters of LDA with a search algorithm would produce better latent topics across words in a document [32].

V. CONCLUSION

In this paper, we gave an overview of web page classification system. Different application areas and an in-depth analysis of the web page classification process were looked into. Analysis of state-of-art techniques for feature selection techniques used in WPC was looked in to with a view to identify challenges fraught by each one. Also related works in the areas of WPC was reviewed to identify the latest works in this domain. It clearly shows that document representation phase is one of the areas that are receiving interest by researchers. Most currently used methods of document representation are Vector Space Model (VSM), Probabilistic Topic Model and Statistical Language Models and Neural network language models [25]. The chosen document representation technique have a direct impact on the classification results as it captures the semantics of document and also contribute to reduce the problem of high dimensionality. Combining different DR technique are new areas of research because each technique perform differently depending on the dataset. Future work in WPC should focus on improving the semantic relationship of web document by hybridizing difference DR technique which will inherently improve the classification result. Also, ontology based techniques can be used to capture the real semantics of web documents.

REFERENCES

- [1] Raj, A. J., Francis, F. S., & Benadit, P. J. (2016). “*Optimal Web Page Classification Technique Based on Informative Content Extraction and FA-NBC*”. Computer Science and Engineering, 6(1), 7-13.
- [2] Deri, L., Martinelli, M., Sartiano, D., & Sideri, L. (2015, November). “*Large scale web-content classification. In Knowledge Discovery*”, Knowledge Engineering and Knowledge Management (IC3K), 2015 7th International Joint Conference on (Vol. 1, pp. 545-554). IEEE.
- [3] Dey Sarkar, S., Goswami, S., Agarwal, A., & Aktar, J. (2014). “*A Novel Feature Selection Technique for Text Classification Using Naïve Bayes*”. International Scholarly Research Notices, 2014.
- [4] Shibu, S., Vishwakarma, A., & Bhargava, N. (2010). “*A combination approach for web page Classification using Page Rank and Feature Selection Technique*”. International Journal of Computer Theory and Engineering, 2(6), 897.
- [5] Dixit, S., & Gupta, R. K. (2015). “*Layered Approach to Classify Web Pages using Firefly Feature Selection by*

- Support Vector Machine (SVM)*. International Journal of u- and e-Service, Science and Technology, 8(5), 355-364.
- [6] Qi, X., & Davison, B. D. (2009). "Web page classification: Features and algorithms". ACM computing surveys (CSUR), 41(2), 1-31.
- [7] Abdelbadie B., Abdellah I., & Mohammed B. (2013). "Web Classification Approach Using Reduced Vector Representation Model Based On Html Tags". Journal of Theoretical and Applied Information Technology, 55(1).
- [8] AbdulHussien, A. A. (2017). "Comparison of Machine Learning Algorithms to Classify Web Pages". International Journal of Advanced Computer Science and Applications (IJACSA), 8(11),
- [9] Mangai, J. A., Kothari, D. D., & Kumar, V. S. (2012). "A Novel Approach for Automatic Web Page Classification using Feature Intervals". International Journal of Computer Science Issues (IJCSI), 9(5).
- [10] Huang, C. C., Chuang, S. L., & Chien, L. F. (2004). "Using a web-based categorization approach to generate thematic metadata from texts". ACM Transactions on Asian Language Information Processing (TALIP), 3(3), 190-212.
- [11] Mangai, J. A., Kumar, V. S., & Balamurugan, S. A. (2012). "A novel feature selection framework for automatic web page classification". International Journal of Automation and Computing, 9(4), 442-448.
- [12] Haveliwala, T. H. (2003). "Topic-sensitive pagerank: A context-sensitive ranking algorithm for web search". IEEE transactions on knowledge and data engineering, 15(4), 784-796.
- [13] Cui, H., Kan, M. Y., Chua, T. S., & Xiao, J. (2004, July). "A comparative study on sentence retrieval for definitional question answering". In SIGIR Workshop on Information Retrieval for Question Answering (IR4QA) (pp. 383-390).
- [14] Hammami, M., Chahir, Y., & Chen, L. (2003, October). "WebGuard: Web based adult content detection and filtering system". In Web Intelligence, 2003. WI 2003. Proceedings. IEEE/WIC International Conference on (pp. 574-578). IEEE.
- [15] Kato, R., & Goto, H. (2016, March). "Categorization of web news documents using word2vec and deep learning". In Proceedings of the 2016 International Conference on Industrial Engineering and Operations Management Kuala Lumpur, Malaysia.
- [16] Fatima, S., & Srinivasu, B. (2017). "Text Document categorization using support vector machine".
- [17] Wang, Z., Ma, L., & Zhang, Y. (2016, June). "A Hybrid Document Feature Extraction Method Using Latent Dirichlet Allocation and Word2Vec". In Data Science in Cyberspace (DSC), IEEE International Conference on (pp. 98-103). IEEE.
- [18] Khan, A., Baharudin, B., Lee, L. H., & Khan, K. (2010). "A review of machine learning algorithms for text-documents classification". Journal of advances in information technology, 1(1), 4-20.
- [19] Alamelu Mangai, J., Santhosh Kumar, V., & Sugumaran, V. (2010). "Recent Research in Web Page Classification—A Review". International Journal of Computer Engineering & Technology (IJCET), 1(1), 112-122.
- [20] Yin, D., Hu, Y., Tang, J., Daly, T., Zhou, M., Ouyang, H., & Langlois, J. M. (2016, August). "Ranking relevance in yahoo search". In Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (pp. 323-332). ACM.
- [21] Socher, Richard, Perelygin, Alex, Wu, Jean Y., Chuang, Jason, Manning, Christopher D., Ng, Andrew Y., and Potts, Christopher (2013). "Recursive deep models for semantic compositionality over a sentiment treebank". In Conference on Empirical Methods in Natural Language Processing.
- [22] Maas, A. L., Daly, R. E., Pham, P. T., Huang, D., Ng, A. Y., & Potts, C. (2011, June). "Learning word vectors for sentiment analysis". In Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies-Volume 1 (pp. 142-150). Association for Computational Linguistics.
- [23] Biro, I., Benczur, A., Szabo, J., & Maguitman, A. (2008, October). "A comparative analysis of latent variable models for web page classification". In Latin American Web Conference, 2008. LA-WEB'08. (pp. 23-28). IEEE.
- [24] Lilleberg, J., Zhu, Y., & Zhang, Y. (2015, July). "Support vector machines and word2vec for text classification with semantic features". In Cognitive Informatics & Cognitive Computing (ICCI* CC), 2015 IEEE 14th International Conference on (pp. 136-140). IEEE.
- [25] Singh, K. N., Devi, H. M., & Mahanta, A. K. (2017). "Document representation techniques and their effect on the document Clustering and Classification: A Review". International Journal of Advanced Research in Computer Science, 8(5).
- [26] Jindal, R., Malhotra, R., & Jain, A. (2015). "Techniques for text classification: Literature review and current trends". Webology, 12(2), 1.
- [27] Azam, N., & Yao, J. (2012). "Comparison of term frequency and document frequency based feature selection metrics in text categorization". Expert Systems with Applications, 39(5), 4760-4768.
- [28] Deerwester, S., Dumais, S. T., Furnas, G. W., Landauer, T. K., & Harshman, R. (1990). "Indexing by latent semantic analysis". Journal of the American society for information science, 41(6), 391.
- [29] Hofmann, T. (1999, August). "Probabilistic latent semantic indexing". In ACM SIGIR Forum (Vol. 51, No. 2, pp. 211-218). ACM
- [30] Le, Q., & Mikolov, T. (2014, January). "Distributed representations of sentences and documents". In International conference on machine learning (pp. 1188-1196).
- [31] Blei, D. M., Ng, A. Y., & Jordan, M. I. (2003). "Latent dirichlet allocation". Journal of machine Learning research, 3(Jan), 993-1022.
- [32] Dit, B., Panichella, A., Moritz, E., Oliveto, R., Di Penta, M., Poshyvanyk, D., & De Lucia, A. (2013, May). "Configuring topic models for software engineering tasks in tracelab". In Traceability in Emerging Forms of Software Engineering (TEFSE), 2013 International Workshop on (pp. 105-109). IEEE.
- [33] Mikolov, T., Chen, K., Corrado, G., & Dean, J. (2013). "Efficient estimation of word representations in vector space". arXiv preprint arXiv:1301.3781.
- [34] Moiseev, G. (2016). Classification of E-commerce Websites by Product Categories. In AIST (Supplement) (pp. 237-247).
- [35] Turney, P. D., & Pantel, P. (2010). "From frequency to meaning: Vector space models of semantics". Journal of artificial intelligence research, 37, 141-188.