

# New technique for Web page Information Categorization using Unsupervised Clustering

Neeraj Mehta, Avinash Rathore  
IES-IPS Academy, Indore

**Abstract :** Classification of web content is dissimilar in a number of characteristic as compared with web page classification. The unrestrained nature of web content nearby added challenge to web page classification as compared to traditional text classification. The web content is semistructured and encloses arrange information in form of HTML tags. A web page consists of hyperlinks to position to other pages. This consistent environment of web pages provide features that can be of superior facilitate in classification. primary all HTML tags are removed from the web pages, together with punctuation marks. The subsequently step is to remove stop words as they are frequent to every documents and does not give a lot in searching. In nearly all cases a Firefly Web Page Classification is functional to diminish words to their basic stem. One such frequently used stemmer is the Firefly Web Page Classification. We proposed Web page Information Categorization (WPIC) algorithm can similarly attain the Categorization of information for dissimilar content .Merge WPIC with personalized search engine technology, and improving the efficiency of WPIC.

Keywords: WPIC , URL, text classification, Firefly Web Page Classification.

## I. INTRODUCTION

Currently, the web pages are increasing at an exponential rate and can cover nearly any information desirable. though, the huge amount of web pages construct it further and more complicated to successfully discover the goal information for a user. usually two solution subsist, hierarchical browsing and keyword searching. though, these pages be different to a enormous extent in both the information content and quality. Furthermore, the association of these pages does not agree to for simple search. So an resourceful and precise technique for classifying this huge quantity of data is extremely necessary if the web pages is to be exploited to its occupied potential. This has been felt for a long time and numerous technique have been tried to resolve this problem. various dissimilar Machine learning based algorithms have been functional to the webpage information classification task, including k-Nearest Neighbour (k-NN Algorithm) [2], , Neural Networks [3], and decision trees [4] Bayesian algorithm [5], Support Vector Machine (SVM) [6]. identify the end user objective in Classical technique web of web page document classification are not suitable for web document classification. lots of of documents on the Web are to small or suffer from a require of linguistic data. This work treat with this problem in two novel technique experiment have prove that hypertext links in web documents frequently direct to documents with comparable semantic content. This study leads to use these referenced

web pages as an extension of the investigated one for the purpose of processing their linguistic data as well. though there are a number of restrictions. The referenced documents have to be placed on the same server and a level of recursion must be limited.

The previous technique increases quantity of linguistic data for the nearly all part of documents sufficient but there is another problem. To use machine learning Algorithms we require to construct a high dimensional vector space where every dimension represent one word from or phrase. In spite of the information that several machine learning algorithms are familiar to elevated number of dimensions, in this case the high number of dimensions decrease algorithm accuracy and Informational - The objective of the user is to get together a number of information from one or added web pages.

Transactional - The intent is to achieve some web- mediate activity, like downloading les, purchase

Items online, etc. base on the additional than taxonomy, primary proposed an automatic query goal classification scheme to distinguish merely amongst navigational and informational queries. in classifying queries among navigational and informational classes by allowing for click distribution and anchor link distribution for automatic query classification. Automatic Web page Categorization significant functions in Internet information and Categorization exploration. This work we will discover out the three subsequent works. The primary the design and implementation of algorithm- Web page Information Categorization (WPIC). Training and classification are two critical stages of WPIC. The subsequent is the application of WPIC in an E-government private organization system. The research to improve the accuracy rate of WPIC and extend the application of WPIC in information systems more than E-governments, such as E-commerce systems. Somewhat improve the algorithm can equally attain the Categorization of information for dissimilar content .Merge WPIC with personalized search engine technology, and improving the efficiency of WPIC.

## II. RELATED WORK

Eda Baykan in at al[1] compared the performances of the dissimilar URL-based language classifiers along a variety of dimensions such as features, algorithms, and training size. as well tested our greatest performing classifiers on ODP + SER dataset on the classification of multimedia Web pages and in small-scale language-focused crawlers. Summarize our major consequences for URL-based Web page language classification.

Ajay S. Patil in at al[2] This paper illustrate Naïve Bayesian (NB) technique for the automatic classification of web sites base on content of home pages. The NB technique, is one of the nearly all effective and straightforward techniques for text document classification and has exhibit superior consequences in preceding study conduct for data mining.

Ariyam Das in at al[3] classified user purpose into three classes with superior precision based on the narration of how users respond to prior search consequences. As the consequence illustrate, majority of the queries issue to a search engine have a conventional unambiguous goal which can be recognized to a great extent by our classified. Many researcher proposed different technique near a multi-cost-sensitive learning for visual quality categorization and a multi-value regression for visual excellence score obligation. Our experiment evaluate the extract features and terminate that the Web page's explain visual features (LV) and text visual features (TV) are the most important disturbing factors toward Web page's visual excellence. Supervised and unsupervised learning to categorize queries as informational, not informational, or ambiguous. the majority of these technique have not measured every classes jointly to recognize the user goal in web queries. primary measured all the three classes but did not look clear of the query and url for the classification purpose. We construct ahead the final work based on the perception that the consumer goal for a specified query may be educated from how users in the past have interact with the returned consequences for the query. The thought following this technique is to discover a distribution over the observed features which explain the observed data but which as well try to maximize the uncertainty, in this distribution. These results in a constrained optimization problem which is then solved using an Web page Information Categorization algorithm.

### III. PROPOSED METHODOLOGY

In this paper we proposed Web page classification technique. Still although content based topic classifiers gave enhanced consequences than URL-based ones, concern classification from URL is preferable when the content is not available, or when classification has the major significance. We can findings for URL-based Web page topic classification as follow. We illustrate that dictionary-based , Firefly Web Page Classification algorithm are not sufficient for higher performance URL-based page classification. For the dictionary-based technique still the best-performing alternative using On the other hand for topic classifiers where precision is important and a measure of recall can be sacrifice, token-based statistical dictionaries can be used. We demonstrate that the features have additional impact on the classifier arrangement than the classification algorithms. Resulting from URLs was the nearly all excellent feature set, considerably better than token. We report a performance which recover the nearly all outstanding formerly reported URL-only management for a diminutive dataset of shopping website web pages by with outline of Web pages for training and testing lead to a huge improvement over

with merely URLs, On the other hand the performance of URL based web page classification decrease when the summary of Web pages are use in training phase in accumulation to URLs. The motivation for this is the dissimilarity among URLs and the outline of Web pages.

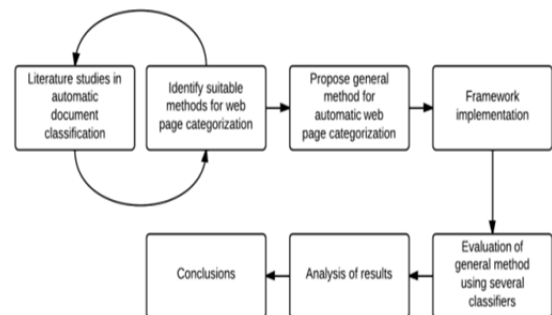


Figure 1: categorization using unsupervised clustering.

We achieve an additional small improvement with using in link information.

Algorithm WPIC in Classification (set of keywords)

Primary Feature Selection (keywords)

Choose the features using term frequency technique return initial set of features; subsequent Feature collection (initial set of features)

Minimize the number of features by selecting the most relevant ones by using evaluators and search technique return concluding set of features;

Classification (concluding set of features)

Classify with semisupervise classifiers Maximize classification accuracy with minimal number of features.

apply boosting to combine dissimilar algorithms gave a small performance improvement of in F-measure. The dispute for URL-based topic classification is. data stability as the definitions of topic be different from one dataset to another, overlap among dissimilar topics in one dataset, empty URLs consisting of only stop tokens or previously.

We primary the generally features that can be used to differentiate navigational, informational and transactional pages. Url Features These are mostly used to identify navigational pages. Navigational pages, being homepages of websites, generally have distinguishing url features such as smaller url depth, url length, occurrence of query keyword in the domain name, etc. HTML Features Different html elements such as tables, images, download buttons, etc. dominate in trans- actional pages. These html features along with the presence of other prominent features can help in different ferentiating transactional pages. Lexical Features Features such as words and sentences per paragraph, amount of text per paragraph, etc. dominate in informational pages. These along with html features are particularly helpful in distinguishing between transactional and informational pages. Bag of Words Features - Transactional and informational query classes can be tied with some specification Keywords. These words are manually selected and weighted differently depending on its occurrence in meta text, title text, headings, special text, anchor text, alternate text and input text. Studied extensively on how to utilize these bag of words features to identify transactional pages. For example,

frequent occurrences of keywords such as buy, cart, on-line store, etc. can indicate a transactional page. Few words like homepage, welcome, etc. can be used to distinguish navigational pages. However, one cannot rely on the bag of words features to identify informational pages as frequently occurring keywords for all domains of information are not easily predictable.

**RESULT ANALYSIS**

After comparison between the Firefly web page classification and in our Method Web page information categorization we have obtained results algorithm by executing a web application. We build this method in two different environments. First is Firefly web page classification and second is Web page information categorization (WPIC). After comparing we found the following advantages over the Firefly web page algorithm.

- i. WPIC programmability greatly reduces testing time for categorization of web.
- ii. Index time calculation is reduced with the help of WPIC over another algorithm.
- iii. Time for data fetching from web page is comparatively less in WPIC.

**Points for Analysis:**

- 1. Based on total time for testing performances.
- 2. Speed Index has better result in WPIC.

**Speed Index:**

The Speed Index is the mean time at which perceptible parts of the page are exhibit. It is indicated in milliseconds and relay on size of the view port.

The Speed Index metric was added to Webpage test in April, 2012 and compute how fast the page contents are visually populated (where lower numbers are better). It is certain convenient for juxtapose experiences of pages against each other (before/after optimizing, my site Vs competitor, etc) and should be used in amalgamation with the other metrics (load time, start render, etc) to better understand a site's performance.

The speed index takes the ocular progress of the perceptible page loading and computes an all inclusive score for how rapidly the content painted. To do this, first it requires to be able to calculate how "complete" page is at diverse points in time throughout the page load.

**FIREFLY CLASSIFICATION**

Following screenshots represent Firefly work implementation:

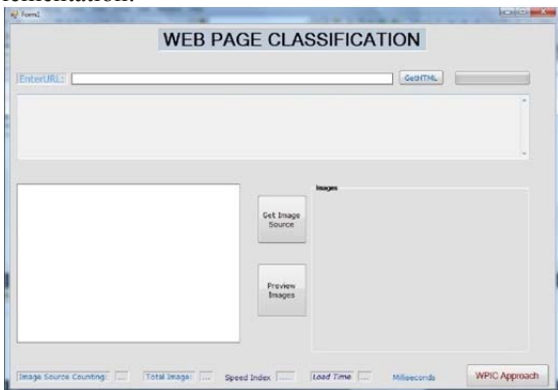


Fig. 2: Firefly Web Page Categorization

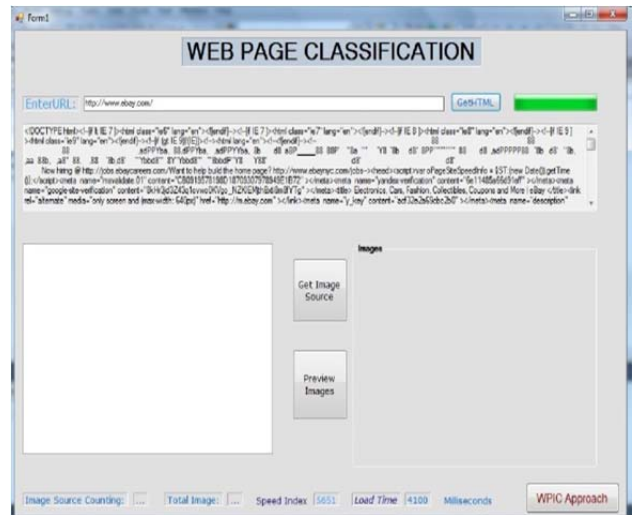


Fig 3: HTML Source Code in Firefly Web Page Categorization Implementation

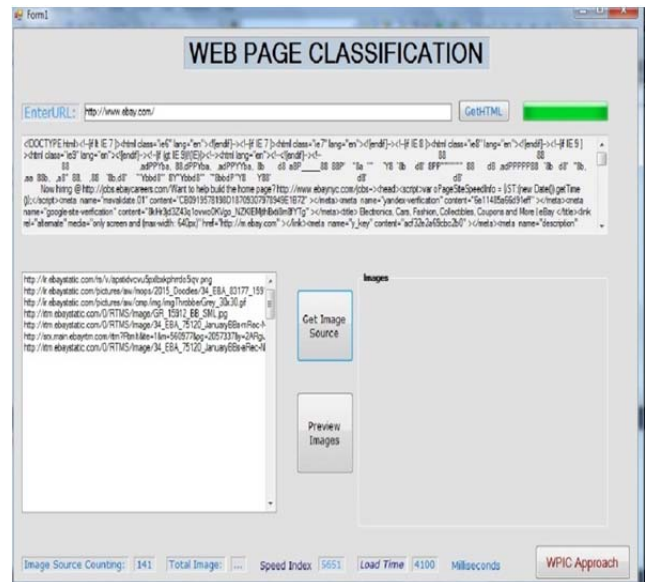


Fig 4: Source and Url in Firefly Web Page Categorization

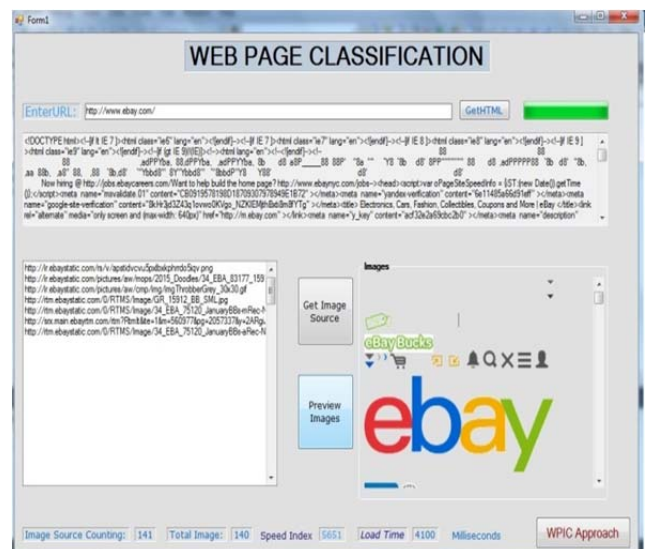


Fig 5: Images on firefly Web Page Categorization

In above diagram we show how firefly classification works in web page classification. In figure 2 shows simple firefly classification tool where one can input web URL and can get HTML source code.

This HTML code extraction show in figure 3. Where we can see that after get HTML process we can see HTML code for the web page. Now firefly use this code to classify image source and image data source which is forwardly shows in figure 4 and 5 so we can get all image sources and can also preview these images.

At bottom of firefly classification tool we can see details about load time, speed index, total Images and image source counting.

From this classification we show an example of www.ebay.com website. Where firefly classification takes 4100 milliseconds load time for classify ebay.com with an speed index of 5651 in loading 140 images.

**WPIC CLASSIFICATION**

Following experimental results show comparison between firefly classification and WPIC classification.

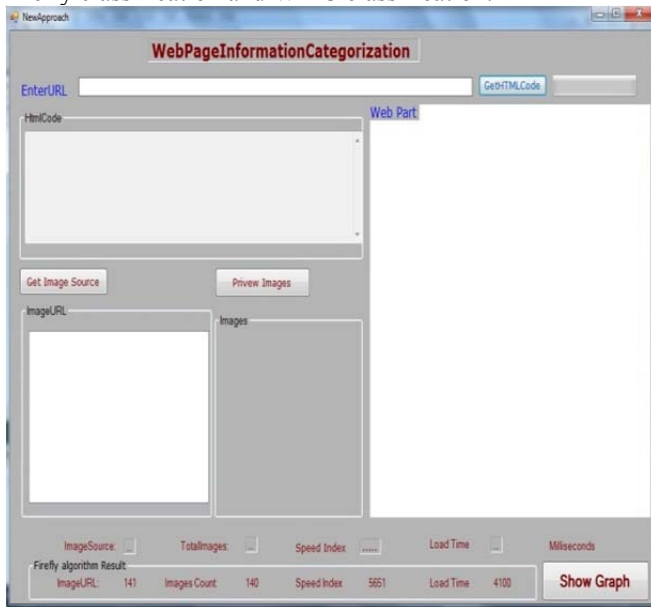


Fig 6: our method (WPIC Approach) Implementation.

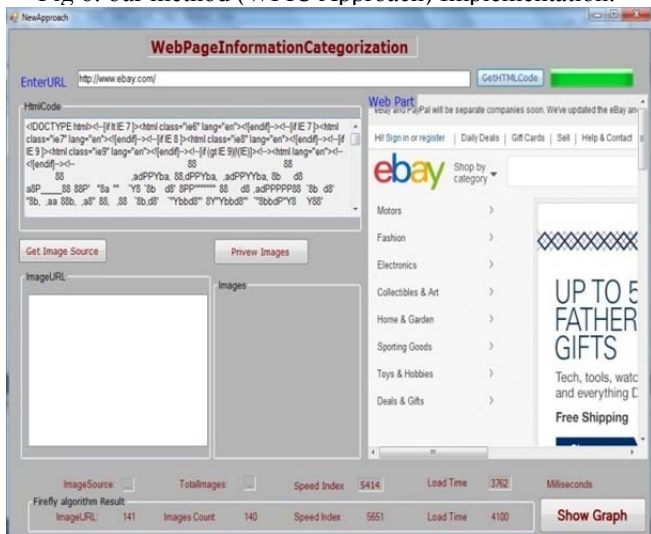


Fig 7: HTML Source Code in WPIC Approach

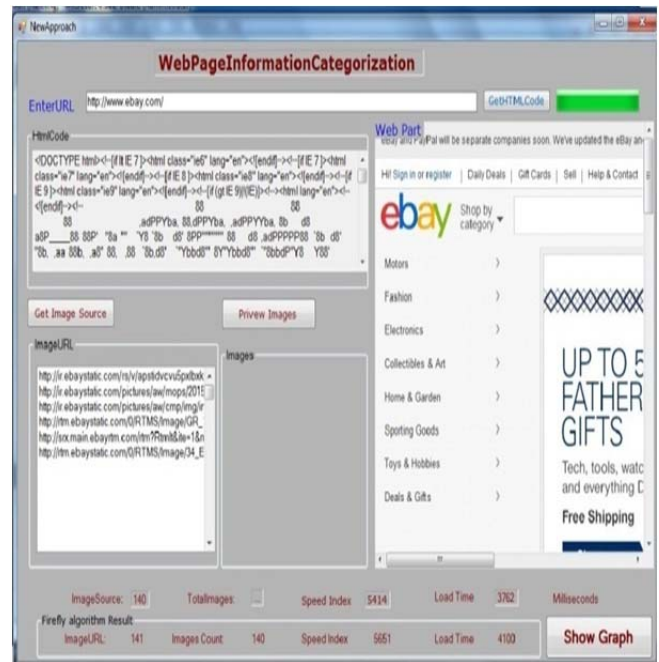


Fig 8: Source and Url in WPIC Approach

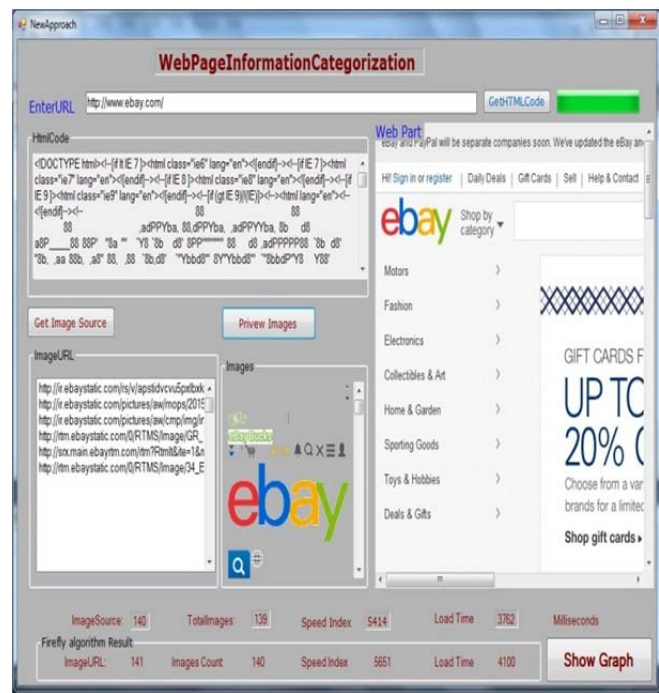


Fig 9: Images on WPIC Approach

As seen in above diagram our developed classification tool can increase speed index with lesser load time.

Figure 6 shows tool for classification where we add extra specification for showing web page also. After getting HTML code we can also see web page with all image sources and images.

In above diagram 9 show comparison between firefly algorithm and WPIC classification where our speed index is 5414 and load time is 3762.

Below table represent comparison between firefly optimization and WPIC classification

Data Source	FIREFLY CLASSIFICATION	WPIC CLASSIFICATION
IMAGE SOURCE	140	140
IMAGE COUNT	141	141
LOAD TIME	4100 ms	3762 ms
SPEED INDEX	5651	5414

Table 1: Comparison between firefly classification Vs WPIC classifications

We also represent our result by graphical analysis where comparison between both classifications can be representing graphically.

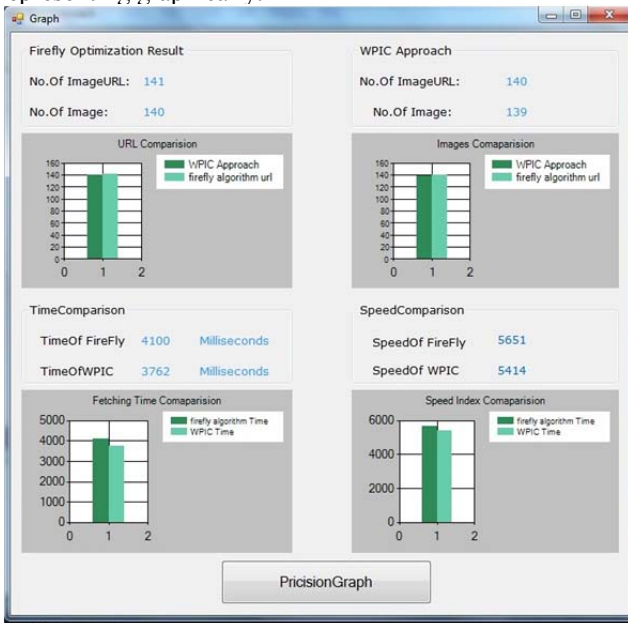


Fig 10: Comparison for Firefly Web Page Classification and WPIC Approach

Above representation shows graphical analysis between firefly classification and WPIC classification where we compare all four parameters URL comparison, Image comparison, Fetching time comparison and speed comparison.

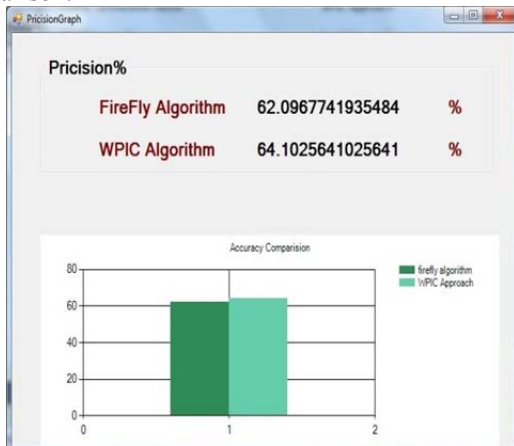


Fig 11: Precision between Firefly Web Page Classification and WPIC Approach

Above graph show precision between firefly classification and WPIC classification. Firefly classification gives 62% accuracy where our WPIC gives 64% accuracy.

CONCLUSION

This technique presented an proficient technique for web page classification. This technique added effective is the training set is set in such a technique that it produce added sets. Though the experimental consequences are relatively encouraging, it would improved if the work with superior data sets with added classes. The existing technique require added or less data for training as well as less computational time of these techniques. After comparing we found the following advantages over the Firefly web page algorithm. WPIC programmability greatly reduces testing time for categorization of web. Index time calculation is reduced with the help of WPIC over another algorithm. Time for data fetching from web page is comparatively less in WPIC.

REFERENCE

- [1.] Eda Baykan, Monika Henzinger, Ingmar Weber” A Comprehensive Study of Techniques for URL-Based Web Page Language Classification” ACM Transactions on the Web, Vol. 7, No. 1, Article 3, Publication date: March 2013.
- [2.] Ajay S. Patil, B.V. Pawar,” Automated Classification of Web Sites using Naive Bayesian Algorithm”IMECS-2012.
- [3.] Ariyam Das, Chittaranjan Mandal, Chris Reade,” Determining the User Intent Behind Web Search Queries by Learning from Past User Interactions with Search Results” The 19th International Conference on Management of Data (COMAD), 19th21st Dec, 2013 at Ahmedabad, India.
- [4.] Ou Wu, Yunfei Chen, Bing Li, Weiming Hu,” Evaluating the Visual Quality of Web Pages Using a Computational Aesthetic Approach” WSDM’11, February 9–12, 2011, Hong Kong, China.
- [5.] Qasem A. Al-Radaideh, Eman Al Nagi “Using Data Mining Techniques to Build a Classification Model for Predicting Employees Performance”, International Journal of Advanced Computer Science and Applications(IJACSA), Vol. 3, No. 2, 2012
- [6.] Thair Nu Phyu “Survey of Classification Techniques in Data Mining”, Proceedings of the International MultiConference of Engineers and Computer Scientists 2009 Vol I IMECS 2009, March 18 - 20, 2009, Hong Kong.
- [7.] W. Chen, Y. Du, P. Zhang, B. Han, "The Effective Classification of the Chinese Web pages based on kNN", JCS, Vol. 6, 2010, pp. 2925-2932.
- [8.] Esra Saraç, Selma Ayşe Özel,” Web Page Classification Using Firefly Optimization” Innovations in Intelligent Systems and Applications (INISTA), 2013 IEEE International Symposium on-Page(s):1 - 5 Print ISBN: 978-1-4799-0659-8.
- [9.] Jon’as Krutil, Milo’s Kud’elka and V’aclav Sn’a’sel,” Web Page Classification based on Schema.org Collection” Computational Aspects of Social Networks (CASoN), 2012 Fourth International Conference on- Page(s): 356 – 360.
- [10.] Selma Ayşe Özel,” A Genetic Algorithm Based Optimal Feature Selection for Web Page Classification” Innovations in Intelligent Systems and Applications (INISTA), 2011 International Symposium on- 10.1109/INISTA.2011.5946076.
- [11.] Win Thanda Aung, Khin Hay Mar Saw Hla,” Random Forest Classifier for Multi-category
- [12.] Classification of Web Pages” Services Computing Conference, 2009. APSCC 2009. IEEE Asia-Pacific- Page(s): 372 – 376.
- [13.] Shiqun Yin Fang Wang Zhong Xie Yuhui Qiu,” Study on Web-page Classification Algorithm Based on Rough Set Theory” Information Processing (ISIP), 2008 International Symposiums on- Page(s):202 – 206.
- [14.] Weitong Huang’ LuXiongXu Junfeng Duan, Yuchang Lu,” Chinese Web-page Classification Study” 2007 IEEE International Conference on Control and Automation ThC2-5 Guangzhou, CHINA - May 30 to June 1, 2007.

- [15.] Suman Roy, A. S. M. Sajeev, Sidharth Bihary and Abhishek Ranjan," An Empirical Study of Error Patterns in Industrial Business Process Models" IEEE TRANSACTIONS ON SERVICE COMPUTING, VOL. XXX, NO. XXX, MONTH 2013.
- [16.] L. W. Han and S. M. Alhashmi, "Joint Web-feature (JFEAT): A Novel Web page Classification Framework", Communications of the IBIMA, Vol. (2010), Article ID 73408, 2010.