

## Proposed Method for Web Pages Clustering using Latent Semantic Analysis

<sup>1</sup>Tawfiq A. Al-Asadi, <sup>2</sup>Ahmed J. Obaid, <sup>2</sup>Ahmad A. Al-Khayatt

<sup>1</sup>College of Information Technology, Babylon University, Babel, Iraq

<sup>2</sup>Department of Computer Science, College of Education, Kufa University, Najaf, Iraq

---

**Abstract:** With explosive growth of information in World Wide Web and large host numbers continuously added to the internet which has different types of content such as text, images, audio, videos and others made the process of analysis web data very complex and difficult process. Analysis of web data also is very important task which required by many organizations, academic centers, companies, agencies, etc., for various task such as enhanced searching process, monitoring and business application. Therefore, many algorithms has proposed in recent research works to construct a model used for knowledge discovering from a corpus of web pages which might exist in different structured forms. The process of extracting web pages features and grouping similar web pages having similar interesting features called web page clustering. In this study, latent semantic analysis has been considered which has applied successfully in many text documents applications by first applying web pages mining process to our selected web pages corpus then text mining process are applied to represent web pages in a Vector Space Model (VSM), finally, K-means algorithm proposed to grouping similar web pages which are exist in semantic space.

**Key words:** Web Content Mining (WCM), web page clustering, Latent Semantic Analysis (LSA), K-means, mining, successfully

---

### INTRODUCTION

Web pages clustering process can be considered one of the most complex and important task in web mining which require to types of data to be studied and analyzed successfully to enhanced services provided on the web. Web users made large number of requests to various web pages which either belong to single web site or different web sites hosted in many web servers on the world. Structure and appearance of web sites contents are differ depend on the design technique used where many web sites designed by using CMS scripts or other languages such as Asp, PHP, HTML, etc. The application of data mining in web data is called web mining. Web data Mining classified into three categories are: Web Usage Mining (WUM), Web Content Mining (WCM) and Web Structure Mining (WSM) based on data type used in mining process. Web usage mining is the application of data mining to understanding user's behaviors and their interesting towards various web sites contents. Web content mining is the task of applying data mining techniques to extracting useful information from the contents of web documents while web structure mining is the process of extracting knowledge about hyperlinks among web pages and the travers paths of users among large number of web pages which may belong to same web site or others sites (Cooley *et al.*, 1997; Al-Asdi and Obaid, 2016).

A starting point to perform web pages clustering process is to create a model for representing web pages by its features, one of the most usage model for representing text documents as well as web pages is VSM where web pages are parsed to extract unique words from set of web pages which has considered as unique features that has a discriminative power to recognize pages among each other's. The most usage weighting schema which is commonly used to weighting features and its important in web pages is by use term frequency inverse document frequency approach (Anastasiu *et al.*, 2013). When all words that extracted from web corpus has been considered, then the VSM become very sparse and most values in VSM are 0's, this phenomena called high-dimensional sparsity (Curse of Dimensionality) and this issue made data mining techniques such as clustering a very difficult process (Charu, 2015a, b). Therefore, most of mining results doesn't have a meaning with present lot of noise features in high-dimensional space. Feature selection/transformation are a pre-processing methods used to reduce high-dimensional space into lower dimension which lead to improve the performance of data mining techniques such as clustering (Shah and Mahajan, 2012). However, reduction of dimension space not only eliminate the noise (irrelevant features) and redundant attributes but also improve the computation space and storage requirement (Alelyani *et al.*, 2014; Aggarwal and Zhai, 2012).

The most commonly method used in dimension reduction where new features are linear combination to the features in original space is Latent Semantic Analysis (LSA). LSA is widely used in document clustering, classification, indexing and information retrieval and is a quite appropriate for a sparse data as in text documents. LSA which is also called latent semantic indexing (when used in information retrieval) employed by search engines to retrieve the relevant result according to user query (Charu, 2015a, b). In this study, web sites of KUFA University has been considered to be processed, the study proposed start crawling the web pages corpus from university web server to local repository then applying pre-processing and analyzing steps where the textual features has been extracted from web pages after pre-processed it then using LSA to extract and represent corpus in semantic space. Finally, point assignment clustering algorithm (K-means) proposed to partition similar group of web pages depend on its correlation to semantic concepts.

**Problem statement:** Web content mining techniques has a number of research issues regarding to extract and discovering knowledge from web pages contents that comes in different forms, web clustering currently considered most critical problems in web application and it's highly important for categories, organize, discovering and extracting knowledge from dynamic web resource contents. There are practical challenges in clustering process for web sites clustering a new algorithm needed to combine different type of data forms and preprocessing varies web resources (Yadav and Mittal, 2013; Vakali *et al.*, 2007). There is another kind of challenges related to representation of extracted features from large collection of web pages. Extracted textual feature from web pages faces several additional problems related to high dimensional feature representation and variant length of web pages. Unlike traditional data mining techniques, clustering of web pages required to use sophisticated algorithm doesn't use in any other application and can deal with sparse features which may fail to group interesting content when use another type of algorithms (Anastasiu *et al.*, 2013). Therefore, there is always a persistent need to analyses web pages from time to time to discover similar web pages that belongs to similar or different web sites which can give a knowledge about what are web pages speaking about.

**Literature review:** The explosive growth of web pages in the world make the indexing process of pages by search engines is difficult in order to represent requested knowledge from users in a proper way, also to group similar web pages that share a meaningful data. Therefore,

clustering process is an unsupervised technique used to group similar objects that are similar to each other within similar group and are varied from those in other corpuses or groups. Thus, many researchers are working from this points by Jain *et al.* (1999) provide an spacious survey of various techniques used for web document clustering process. In 2004, Hammouda and Kamel (2004) proposed a clustering algorithm based on Document Index Graph (DIG). Andrews and Fox (2007) illustrate a development approaches used for web document clustering while (Hung and Xiaotie, 2007) measure similarity among web documents by using Suffix-tree similarity measure then used Hierarchal Clustering Algorithm (HCA) to group similar web pages. Srinivas and Mohan (2010) proposed a clustering algorithm by using incremental hierarchical clustering approach. Jensi and Jiji (2014) illustrate an optimization approaches for clustering text documents including LSA concepts which used for indexing similar documents. Nirkhi and Hande (2008) discuss an overview of most used web document clustering algorithms and its application and limitation exist in web environment. Xiao (2010) presented an overview of semantic techniques used in web document clustering and its limitation when applying in real time web data. In this study, we employee LSA concept to extracting a knowledge from real time web pages that collected from Kufa University web sites which has a different sources and content and illustrate the result of applying semantic analysis on our selected data.

## MATERIALS AND METHODS

**Proposed models:** Implementation of our proposed system consist from several steps including crawling web data till clustering analysis model, Fig. 1 shows our proposed model.

**Crawling data:** Web crawler system also called web spider software are the programs can download web pages, data on the web is sparse on million on pages serve by hug number of servers when users follow the hyperlinks to access information from one page to other. Crawler can visit many web sites to collect web data and used for mining and storage purpose. In our proposed model, we use crawling software for catching website pages and images and store it in repository. There are many crawling software that used for different purpose such as gathering information and links updating for web site content, monitoring and performance test and other tasks. One of the most usage crawling software Teleport Pro., it's possible to apply configuration of it for catching the web pages and images and avoiding other files such as CSS files, java files and so on. In images it's also

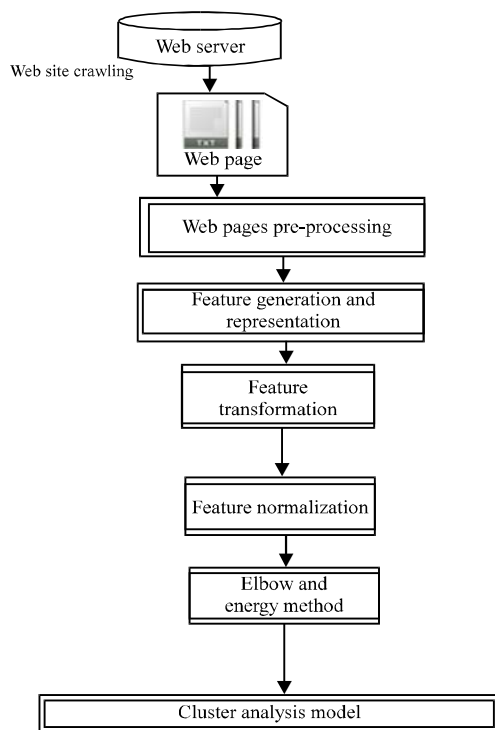


Fig. 1: Proposed model diagram

possible also catching images and filtering it by avoiding such as gif and small size images. The main feature in this software is able to catch any page or content for that web site and store it in static HTML page, some web sites design by use CMS (Content Management Systems), so, articles are store also as a static HTML pages for this reason, we use it for our web site.

**Web pages preprocessing:** Web page like any text document need to be preprocessed before applying mining algorithms such as extracting the terms and counting words in it, the difference in web document such as HTML page is by tag removal and identified the main content blocks in that page and this required carefully and nested procedures when applying preprocessing on HTML pages (Srivastava *et al.*, 2005).

**Preprocessing steps:** Web pages preprocessing include some preprocessing steps as follow.

**Identifying text fields:** There are different text fields in HTML pages such as title, body and metadata. In many search engine search it consider the title of article or keywords when searching about related content in our model we assign weight for title let say 0.5 for the words that appear in title and remaining for words that in body.

**Removing HTML tags:** By removing HTML tags may cause problems by joining text that separated by tags, there are many software can be used for converting HTML pages to text files and then preprocessing text files for extracting text features from it. We introduce this problems in HTML pages and solve it by used new tools for extracting HTML text this tool is called “Html AgilityPack”. By using this tools text in HTML tags can be extracted then store it in file that carried its own name.

**Anchor text identification:** Anchor text is the text that associate with hyperlinks, anchor text contain information that pointed to other pages. In many HTML pages we found many pages contain same anchor text that linkage pages among others for this reason these words were eliminated in text files.

**Main content blocks text identification:** Commercial web pages contain many information regarding to description, copy right and search notes to other related topics. This issue happen with a commercial and different content pages, when using nested pages design and static content, blocks of text content in these HTML pages will be simply identified.

**Text preprocessing:** To preprocess text files need to do again many processes due web pages can be varied and contain many contents some of these pages may be news articles, blogs, information about particular topics, lectures, notes and others while some pages might written in many languages. Clustering or classification of these pages by topics could be a challenges problem. There are many approaches for dealing with document indexing and information retrieval such as Vector Space Model (VSM) and Latent Semantic Indexing (LSA). A practical method for web document representation is by used VSM representation, in this study, we discussed briefly about how to preprocess text of web documents as shown in the following steps (Liu, 2011).

**Stop word removal:** Stop words are commonly occurred words in documents, these words provide little information about documents and it’s actual content, typically stop words commonly available around 300-400 words and its available by Yadav and Mittal (2013).

**Stemming:** Variation of the same word need to be consider when preprocess text files, for example, the variation of the word “Computer” can be found as “Computers” for English languages and for Arabic can be found as “الحاسوب”, “اله حاسبة”, “حاسب الي”, “كو مبيوتر” and others all of these words refer to one general word

category “ computer”. In the model we consider status of two languages due the web site for KUFA University has been designed by using two languages Arabic and English. Articles has also been written by those languages. The main drawback in stemming is in some cases there are many words refer to different meaning, this challenge face us if we try to clustering or classification text document based on its content.

**Punctuation marks:** After stemming process has been done, punctuation marks such as comma, semicolons, numbers, dates, hyphens and special characters are eliminated. Typically hyphenated word can be treated as separated word or they merged into single word. Many common used words may added to dictionary that used for feature extraction process. Result of this step are text documents that contain distinct and meaningful words and will be treated as a bog-of-words which the ordering of the words irrelevant. Use bag-of-words representation is suddenly effective way for fast document labeling and classification.

**Web (pages) documents representation:** In this step each web document is represented by vectors where the rows represent web documents and columns are unique words. The entries of web document-word matrix are frequencies of unique word in every web document. A popular weight for word frequency is by use Tf×Idf schema. In this method, the weight of word (term) $t_{ij}$  web document  $d_j$  is the number of times that word (term)  $t_i$  appears in web document  $d_j$  and it's denoted by  $f_{ij}$ . Idf is the inverse document frequency, this can reflect the importance of given words (terms) that appears in web document this schema sometimes used for text classification (Joachims, 1998). The equations for a given schema as follow:

$$tf_{ij} = \frac{f_{ij}}{\text{Max}\{f_{i1}, f_{i2}, f_{i3}, \dots, f_{i|v|}\}} \quad (1)$$

Where:

$tf_{ij}$  = The normalized-term frequency

$|v|$  = The vocabulary size of the collection

The inverse document frequency is given by:

$$idf_i = \frac{\log N}{df_i} \quad (2)$$

Where:

$N$  = Total number of web documents

$df_i$  = The number of web document that contain the term  $t_i$

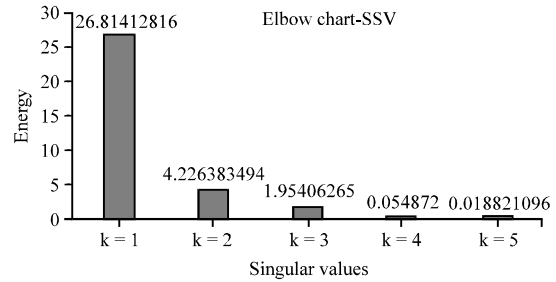


Fig. 2: SSV of word (term) web pages matrix

Then, the weight of term  $t_i$  is given by:

$$w_{ij} = tf_{ij} \times idf \quad (3)$$

If the term  $t_i$  does not appear in web document  $d_j$ , then, otherwise refer to number of occurrence that term  $t_i$  in web document  $d_j$ . In order to obtain better result for web document with different length then the following Eq. 4:

$$w_{ij} = tf_{ij} \times idf = \frac{f_{ij}}{\sum_{s=1}^N (tf_{ij} idf(w_{sj}))^2} \times \log_2 \left( \frac{N}{df_i} \right) \quad (4)$$

Finally, the result as shown in Fig. 2 given to the next step which is called “feature transformation” which is used for reduce the high dimensionality of extracted features.

**Feature transformation:** Feature transformation also called (feature reduction) aim to mapping the features in original space to new feature space has lower dimension and can't linked to the features in the original space. The goal for applying feature transformation to creating intelligent features to avoid the problems in original features such as polysemy, synonymy and homonymy which present in web documents. The most commonly method used in dimension reduction in which the new features are linear combination to the features in original space is Latent Semantic Analysis (LSA). LSA is widely used in document clustering, classification, indexing and information retrieval and is a quite appropriate for a sparse data as in text documents. LSA which is also called latent semantic indexing (used with information retrieval) by search engines to retrieve the relevant result according to user query.

**Normalization and elbow curve plotting:** The result of pre-processing steps and feature transformation using LSA method for input word (term) web pages matrix which is called (A) is three matrices as shown in the following Eq. 5:

$$A = U \sum V^T \quad (5)$$

Where:

U = The mapping of words (terms) to concept view

V<sup>T</sup> = The mapping of web pages to extracted concepts

Therefore, for large collection of data set that contain hundred thousands of web pages (or web documents), recommended dimensions based on recent studies between 100-500 dimensions. Hence, with this number of dimension a new contribution work can be extracted is by applying clustering algorithm on extracted semantic space as in the proposed system where partitioning clustering algorithm proposed to cover the problem of finding similarities among large number of web pages. Elbow curve plotting is a method used in LSA for finding the optimal truncated value in reduction step where K in LSA implementation represent the new truncated space for the matrices of  $\{U \sum V^T\}$  to reduction space of (K)

where  $A'_k = \{U_k \sum_k V_k^T\}$ , elbow method start by plot the Squared of Singular Values (SSV) as shown in Fig. 3 and 4.

In Fig. 3, the first value seems added more significant for information retrieval or analysis content than other values, therefore, the first dimension give an absolute value for the web pages collection. However, this value correspond to the length of web pages in matrix V<sup>T</sup> and correspond to the number of times words used in all web pages in matrix U.

**Cluster analysis model:** One of the most used algorithm especially in text mining is K-means algorithm which used Euclidian distance as default dissimilarity measure. In text mining where number of dimension become very high which is called curse of dimensionality, most feature transformation method used in text mining is LSA which

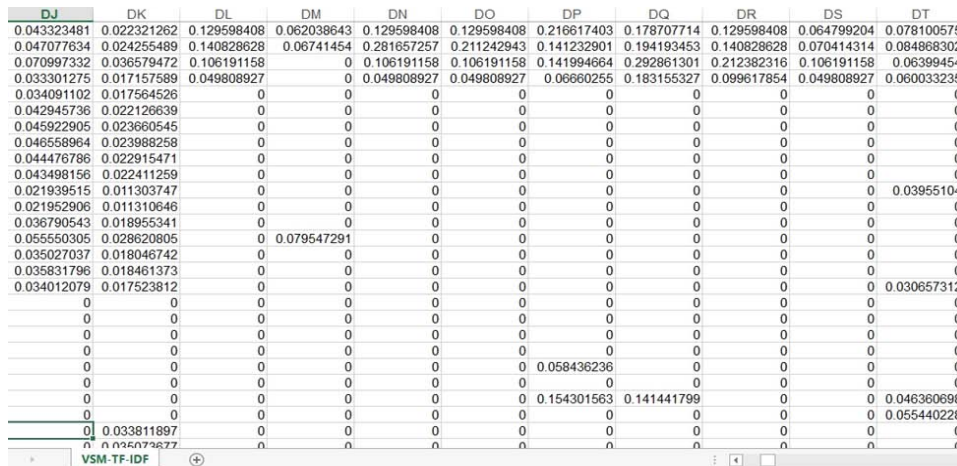


Fig. 3: Web pages in VSM representation

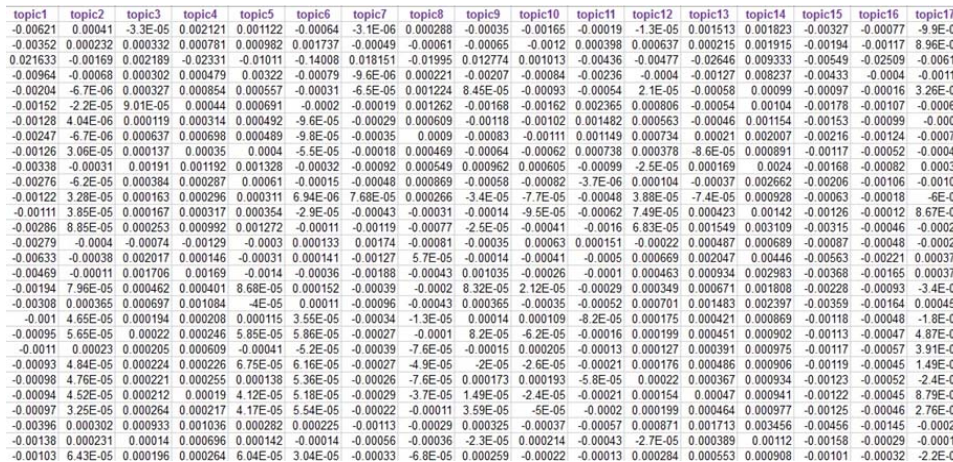


Fig. 4: Matrix (U) words (terms) concept similarity



is used by search engine, information retrieval, text noise reduction and irrelevant feature elimination tasks. The primary goal of K-means algorithm is to partition data set into k-clusters where K here its differ from K used in reduction step, by represent each cluster with representative points that summarize cluster and assign each point to the appropriate cluster (Langhnoja *et al.*, 2013), normalization can give better clustering result due to the various ranges in our data point in each space.

**RESULTS AND DISCUSSION**

After completing pre-processing steps the result of representing web pages corpus is a matrix of web page-words matrix as shown in the following, where the rows representing web pages, unique words (terms) in

columns and entries are weights of every word (term) in web corpus. Clustering process, Fig. 5 represents the clusters for our data set and centroids for each cluster, in these clusters contain some sites which has been browsed by users each cluster represent a different space from others.

After pre-processing steps completed, LSA has implemented which is used a method called Singular Value Decomposition (SVD) on result matrix from Fig. 3, the result of applying LSA produce three matrices as discussed previously which are U,  $\Sigma$  and  $V^T$  as shown in Fig. 6, respectively.

Figure 5 shows the result of singular values representation in descending orders which means the extracted semantic concepts in the beginning have more discriminative power than others, truncated can be done

Fig. 5: Matrix ( $\Sigma$ ) represent singular values in descending order

Fig. 6: Matrix ( $V^T$ ) represent web pages-concepts correlation

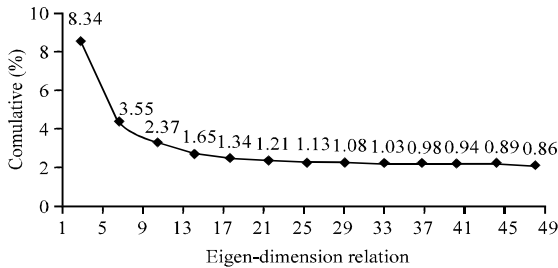


Fig. 7: Elbow Graph of SQR ( $\sigma$ ) dimensions values

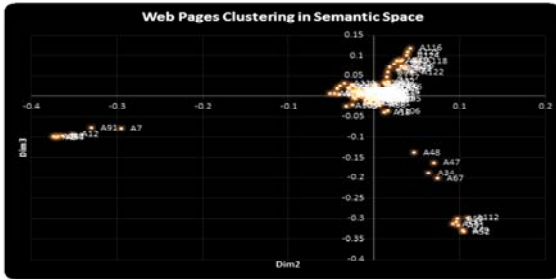


Fig. 8: Web pages clustering in semantic space using K-means

based on the result of this matrix where K has selected based on the number of web pages in the corpus. Figure 6 shows the mapping of web pages in the corpus to new extracted semantic concepts in semantic space, when K has been selected the truncated of this matrix start from bottom of matrix  $V^T$ .

To select appropriate truncated Value for the matrix ( $V^T$ ) which represent the correlation of web pages to the new concepts derived in semantic space, square singular values of matrix ( $\Sigma$ ) used for this purpose which reflect the importance of concepts in descending order that means the highest value assigned to concepts has higher discriminative power than others which might doesn't have that property, Fig. 7 shows the importance of extracted semantic concepts towards 50 dimensions based on singular values importance.

By using reduction of feature space in semantic analysis lot of noise, irrelevant and synonymy issues are eliminated when extracting concept in new concept space. Therefore clustering of web pages in semantic space provide many meaningful results as shown in Fig. 8 when applying partitioning based algorithm such as K-means.

### CONCLUSION

This study deal with the problem of clustering web sites pages based on discovering hidden information from the whole pages itself, clustering of web pages

(documents) are differ from traditional text document clustering which due its contain lot of embedded information in its structure design as well as the extraction of pure text content from these pages is a practical challenges when corpus contain thousands of web pages comes in different design and forms. Therefore, in this study, we introduce a solution for labeling web pages by combination of using semantic analysis of our web pages corpus then applying K-means algorithm which is successfully partitioned our selected data in to set of clusters which given to us what our web pages speaking about and find similar web pages as well as similar web sites specially when labeling property absent from our selected data.

### RECOMMENDATIONS

For a future research, we plan to find group of similar words together by implementing LDA concept, we searching for enhanced topic modeling process by employing non matrix factorization by selecting a proposer truncated values to produce a clusters with real topic labels.

### REFERENCES

Aggarwal, C.C. and C.X. Zhai, 2012. A Survey of Text Clustering Algorithms. In: Mining Text Data, Charu, C.A., and X.Z. Cheng (Eds.). Springer, New York, USA., ISBN:978-1-461 4-3222-7, pp: 77-128.

Al-Asdi, T.A. and A.J. Obaid, 2016. An efficient web usage mining algorithm based on log file data. J. Theor. Appl. Inf. Technol., 92: 215-224.

Alelyani, S., T. Jiliang, L. Huan, 2014. Feature Selection for Clustering: A Review. In: Data Clustering: Algorithms and Application, Charu, C.A. and C.K. Reddy (Eds.). CRC Press, Boca Raton, Florida, ISBN-13:978-1-4665-5822-9, pp: 29-60.

Anastasiu, D.C., A. Tagarelli and G. Karypis, 2013. Document Clustering: The Next Frontier. In: Data Clustering: Algorithms and Applications, Charu, C.A. and K.R. Chandan (Eds.). CRC Press, Boca Raton, Florida, ISBN-13: 978-1-4665-5822-9, pp: 305-338.

Andrews, N.O. and E.A. Fox, 2007. Recent developments in document clustering. MSc Thesis, Virginia Tech, Blacksburg, Virginia.

Charu, C.A., 2015a. Mining Text Data. In: Data Mining the Text Book, Aggaral, C.C. (Ed.). Springer, Switzerland, Europe, ISBN:978-3-319-14141-1, pp: 429-456.

- Charu, C.A., 2015b. Cluster Analysis. In: Data Mining the Text Book, Aggaral, C.C. (Ed.). Springer, Switzerland, Europe, ISBN:978-3-319-14141-1, pp: 153-204.
- Cooley, R., B. Mobasher and J. Srivatsava, 1997. Web mining: Information and pattern discovery on the world wide web. Proceedings of the 9th IEEE International Conference on Tools with Artificial Intelligence, November 3-8, 1997, Newport Beach, CA., pp: 558-567.
- Hammouda, K.M. and M.S. Kamel, 2004. Efficient phrase-based document indexing for web document clustering. IEEE. Trans. Knowl. Data Eng., 16: 1279-1296.
- Hung, C. and D. Xiaotie, 2007. New suffix tree similarity measure for document clustering. Proceedings of the 16th International Conference on World Wide Web, May 08-12, 2007, ACM, Banff, Alberta, ISBN:978-1-59593-654-7, pp: 121-130.
- Jain, A.K., M.N. Murty and P.J. Flynn, 1999. Data clustering: A review. ACM Comput. Surv., 31: 264-323.
- Jensi, R. and D.G.W. Jiji, 2014. A survey on optimization approaches to text document clustering. Intl. J. Comput. Sci. Appl., 3: 31-44.
- Joachims, T., 1998. Text categorization with support vector machines: Learning with many relevant features. Proceedings of the 10th European Conference on Machine Learning, Chemnitz, Germany, April 21-23, 1998, Springer, Berlin, Heidelberg, pp: 137-142.
- Langhnoja, S.G., M.P. Barot and D.B. Mehta, 2013. Web usage mining using association rule mining on clustered data for pattern discovery. Intl. J. Data Min. Tech. Appl., 2: 141-150.
- Liu, B., 2011. Information Retrieval and Web Search. In: Web Data Mining, Liu, B. (Ed.). Springer, Berlin, Germany, ISBN:978-3-642-19459-7, pp: 211-268.
- Nirkhi, S. and K. Hande, 2008. A survey on clustering algorithms for web applications. Proceedings of the 2008 International Conference on Semantic Web & Web Services (SWWS), July 14-17, 2008, CSREA Press, Las Vegas, Nevada, pp: 124-129.
- Shah, N. and S. Mahajan, 2012. Document clustering: A detailed review. Intl. J. Appl. Inf. Syst., 4: 30-38.
- Srinivas, M. and C.K. Mohan, 2010. Efficient clustering approach using incremental and hierarchical clustering methods. Proceedings of the 2010 International Joint Conference on Neural Networks (IJCNN), July 18-23, 2010, IEEE, Barcelona, Spain, ISBN:978-1-4244-6916-1, pp: 1-7.
- Srivastava, T., P. Desikan and V. Kumar, 2005. Web Mining Concepts, Applications and Research Directions. In: Foundations and Advances in Data Mining, Wesley, C. and Y.L. Tsau (Eds.). Springer, Berlin, Germany, ISBN:978-3-540-25057-9, pp: 275-307.
- Vakali, A., G. Pallis and L. Angelis, 2007. Clustering Web Information Sources. In: Web Data Management Practices; Emerging Techniques and Technologies, Vakali, A. and P. George (Eds.). Idea Group Publishing, Hershey, Pennsylvania, USA., ISBN:1-59904-230-4, pp: 34-55.
- Xiao, Y., 2010. A survey of document clustering techniques and comparison of LDA and moVMF. Master Thesis, North Carolina State University, Raleigh, North Carolina.
- Yadav, M. and M.P. Mittal, 2013. Web mining: An introduction. Intl. J. Adv. Res. Comput. Sci. Software Eng., 3: 683-687.