# REVIEW ON WEB MINING TECHNIQUE

[1]**Prageet Bajpai**
Computer Science & Engineering
SSCET Bhilai
Chhattisgarh India

[2]**Siddharth Choubey**
Computer Science & Engineering
SSCET Bhilai
Chhattisgarh India

[3]**Abha Choubey**
Computer Science & Engineering
SSCET, Bhilai
Chhattisgarh, India

*Abstract— Data warehouse is an important contemporary issue for many organizations and is relatively a new field in the realm of information technology. As data warehousing is a new field, little research has been done regarding the characteristics of academic data and the complexity of analyzing such data. Educational institutions measure success very differently from business-oriented organizations and the analyses that are meaningful in such environments pose unique problems in data warehousing. The purpose of this thesis is to provide a security of data ware house In the present work we have introduced a vernam cipher bit wise encryption method.*

## I. INTRODUCTION

As every day utilization of World Wide Web is expanding, mining of the database is having all the more requesting. The database of web is in the type of web sessions with session id or session proprietor. So for this kind of mining is called Web Mining. Also, in the Web Mining when we use to discover way traversal design for choice administration in web composition, at that point it goes under the web utilization Mining . Web utilization mining is likewise called web log mining. It is a web mining strategy which depends on the disclosure and investigation of web utilization designs from web logs. These web logs incorporate web server logs, intermediary server logs, web program logs, and so forth. and are made when clients speak with the web server. Web utilization mining is the procedure of discovering what clients are searching for on the web. Scarcely any clients may be taking a gander at just recorded information, though some others may be keen on sight and sound information. It is the accommodation of statistical data points mining methods to discover fascinating use designs from World Wide Web raw numbers in arrangement to acknowledge and better serve the yearnings of Web based applications .another session in LCS is arranged into any of the bunches and expectation list is produced in view of the route examples of relating bunch. An expectation display by considering request data of pages and time spent on them in a session.Data utilized for Web utilization mining, can be gathered at one of these three sections: Server level, Client Level, furthermore, intermediary level accumulation. The contribution for the Web Usage Mining strategy is a customer session report, which is on a very basic level a pre-prepared record and includes data, for instance, who went to the site and what pages were gone to and for to what degree, with their particular request. Notwithstanding it might contain superfluous data. This superfluous data can be limited or diminished by information pre-preparing of the web log information. After information cleaning the extraordinary number of client and session distinguish and finally phase of information preprocessing we get the successive client's get to design from the web server get to log document data.A specific client or an arrangement of client exploiting learning picked up from client's navigational example and enthusiasm of the individual client with the conjunction of substance and structure of Web Sites.Web utilization mining includes the programmed recognition of client get to designs on at least one web servers. It is an utilization of information mining calculations to web get to logs to discover the patterns and regularities in web clients' route designs. There are numerous sorts of information that can be utilized as a part of web mining furthermore, can be grouped into following five sorts:

In a technique to foresee the client's route designs is proposed utilizing grouping and arrangement from Web log information. In the first place period of this technique concentrates on isolating clients in Web log information, and in the second stage bunching process is utilized to gathering the clients with comparable inclinations. At last in the third stage the consequences of grouping and bunching are utilized to foresee the users" next solicitations.

To beat the downsides of the current recommender framework, for example, knowledge, versatility, adaptability, restriction o f precision, we introduce design for incorporating semantic data about the items with web log information and produce a rundown of suggested items by utilizing LCS and regular example Algorithm. An information cleaning methodology ought to fulfill a few prerequisites. Right off the bat, it should recognize and expel every single significant blunder and irregularity in the database. The approach should upheld by devices to restrict manual review and programming exertion and be extensible to effectively cover extra source.

Moreover, information cleaning ought to perform mapping capacity and consolidating capacity. Mapping capacities for information cleaning and other information changes ought to be determined decisively and be reusable for other information sources and in addition for inquiry preparing. The principle point of utilizing affiliation manage mining procedure is to discover the relationship between things in a certain value-based record. Discovering affiliation rules is completely in view of the support and certainty models, where a base bolster must be indicated to begin the pursuit

## II.BACKGROUND STUDY

Data mining efforts associated with the Web, called Web mining, can be broadly categorized into three areas of interest based on which part of the Web to mine; Web Content mining, Web Structure mining, and Web Usage Mining (Kosala and Blockeel,2000). In Web mining, data can be collected at the server-side, client-side, proxy servers or a consolidated Web/business database(Srivastava et al., 2000). The information provided by the data sources can be used to construct several data abstractions, namely
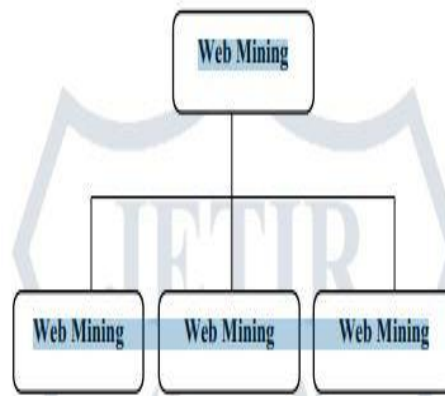
**Fig 1 Web Mining Taxonomy**

users, page-views, click-streams and server Web Usage Mining is defined as the process of applying data mining techniques to the discovery of usage patterns from Web logs data which to identify Web user's behavior (Srivastava et al., 2000). Web Usage Mining is the type of Web mining activity that involves an automatic discovery of user access patterns from one or more Web servers Process of Web Usage Mining: As shown in Fig. 1, three main tasks are performed in Web Usage Mining; Pre-processing,Pattern Discovery and Pattern Analysis. Fig. 1 represents a brief description about the main task of Web Usage Mining process. Web Usage Mining involves determining the frequency of the page access by the clients and then finding the common traversal paths of the users. First task is the data is collected from web server log file.

**Text Mining**

Text mining is a practice that is utilized to find advantageous inarrangemention from large amount of data sets. Data mining has guidelines known as frequent pattern and association rule that is essential for finding frequent patterns. Text Mining is the recognition by computer of new, previously unidentified inarrangemention, by automatically mining inarrangemention from different written resources. Text mining techniques are the fundamental and permitting tools for efficient organization, triangulation, retrieval and summarization of large file quantity. With more and more text, inarrangemention are distribution around on Internet, text mining is rising in importance. Text clustering and text classification are two important tasks in the field of text mining.
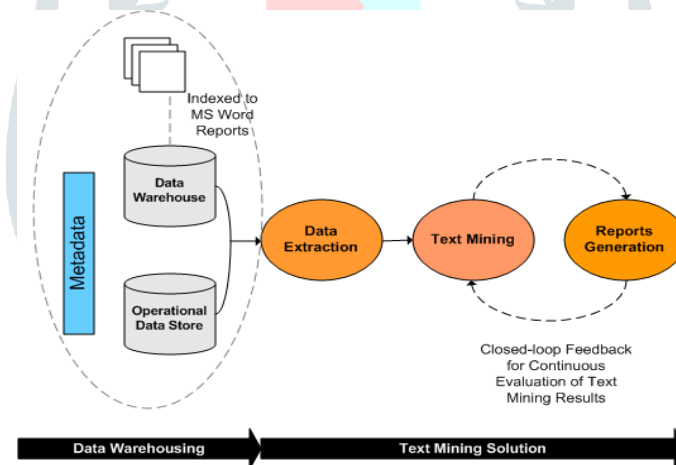


**Fig 2 Process of Text Mining Block Diagram**

Text Clustering is to find out the groups inarrangemention from the text file and cluster these file into the most relevant groups. Text clustering clusters the file in an unsupervised way and there is no label or class inarrangemention. Clustering techniques have to determine the connections between the file and then based on these connections the files are bunched. Given a enormous volumes of files, a superior document clustering techniques may organize those huge statistics of documents into meaningful groups, which permit further browsing and navigation of this quantity be much easier (B.Liu, M.Hu and J. Cheng, 2005). A basic idea of text clustering is to find out which types of documents have many words in common and place these types of documents with the most words in mutual into similar group.

Text Classification is to establish the files into predefined classes with meaningful labels. As text classification wants the facts about those predefined categories, it is applied in a supervised way.

Eighty percent of the inarrangemention in the world is currently stored in amorphous textual arrangement. Although method such as Natural Language Processing (NLP) can complete restricted text analysis, there are currently no computer programs available to investigate and interpret text for the different inarrangemention extraction wants. Thus text mining is a dynamic and unindustrialized region. The world is fast becoming inarrangemention comprehensive, in which specialized inarrangemention is being poised into extremely large data sets. For instance, Internet contains a large amount of online text files, which rapidly change and grow. It is nearly dreadful to manually organize such vast and quickly evolving inarrangemention. The requisite to extract useful and relevant inarrangemention from such bulky data sets has led to a significant requirement to develop computationally competent text mining algorithms (A.M.Popescu and O. Etzioni , 2005). An instance, problem is to automatically dispense natural language text files to predefined sets of categories grounded on their contented. Other instances of problems involving large data sets comprise searching for targeted inarrangemention from technical citation databases (e.g. MEDLINE); search, filter and classify web pages by topic and routing relevant email to the proper addresses.

Text mining is the involuntary and semi-automatic removal of implicit, previously indefinite, and hypothetically useful inarrangemention and patterns, from a large amount of amorphous textual data, such as natural-language text (D. Kerr, H. Mousavi, and M. Iseli ,2013). In text mining, every file is represented as a vector, whose dimension is almost the number of diverse keywords in it, which can be very large. One of the major contests in text mining is to categorise textual data with such superior dimensionality. In adding up to high dimensionality, text-mining algorithms would also deal with word ambiguities such as pronouns, synonyms, and deafening data, spelling mistakes, abbreviations, acronyms and inadequately structured text. Text mining algorithms are of two types: Supervised learning and unsupervised learning. Support vector machines (SVMs) are a set of supervised learning approaches utilized for classification and reversion. Nonnegative matrix factorization is an unsupervised learning method.

## III. LITERATURE REVIEW

Many data mining techniques have been proposed for mining useful patterns in text documents. How to effectively use and update discovered patterns is still an open research issue, especially in the domain of text mining. Since most existing text mining methods adopted term-based approaches, they all suffer from the problems of polysemy and synonymy. Over the years, people have often held the hypothesis that pattern (or phrase)-based approaches should perform better than the term-based ones, but many experiments do not support this hypothesis.

Zhong N. et al (2012) presented an innovative and effective pattern discovery technique which includes the processes of pattern deploying and pattern evolving, to improve the effectiveness of using and updating discovered the patterns for finding relevant and interesting information. Substantial experiments on RCV1 data collection and TREC topics demonstrated that the proposed solution achieves encouraging performance.

Luepol Pipan maekaporn (2013), presented a novel pattern mining approach to RF. This approach mined patterns in both positive and negative feedback and then classified them into clusters to find user-specific patterns. They also proposed a novel pattern deploying method that effectively used the discovered patterns for improving the performance of searching relevant documents. Experiments are conducted on Reuters Corpus Volume 1 data collection (RCV1) and TREC filtering topics. The results shown that the proposed approach achieves promising performance comparing with state-of-the art term-based methods and pattern-based ones.

They also applied a novel pattern deploying the strategy to improve the performance of frequent patterns in text. They evaluated the proposed approach by using it to discover high-quality features in relevance feedback for improving the information filtering. Their results on RCV1 data collection and TREC filtering topics confirmed that the best improvements are obtained by our approach compared to state-of-the-art term-based methods and pattern-based ones.

Bhushan Inje, Ujawla Patil (2014) examined and investigated this fact with considering several states of art data mining methods that gives satisfactory results to improve the effectiveness of the pattern. Here they implemented the pattern detection method to solve problem of term-based methods and improved result which is helpful in information retrieval systems. Their proposal was also evaluated for several they'll distinguish domain, offering in all cases, reliable taxonomies considering precision and recall along with F-measure. For the experiment, they used Reuters (RCV1) dataset and the results show that they improved the discovering pattern as compared to previous text mining methods. The results of the experiment setup show that the keyword-based methods not give better performance than pattern-based method. The results also indicated the removal of meaningless patterns not only reduces the cost of computation but also improved the effectiveness of the system.

In this paper, they have investigated the existing data mining methods with respect to the alternating approach for finding relevant pattern in large documents collection; some research work have been used phrases rather than individual words. However, the effectiveness of the text mining systems was not improved very much. The likely reason is that, a phrase-based method has "lotheyr consistency of assignment and lotheyr document frequency for terms". Hence, in this paper, they presented a concept for mining text documents for sequential patterns. Instead of using single words, they used pattern-based taxonomy (is-a) relation to represent documents. By pruning meaningless (negative) patterns, which have been proven the source of the 'noise' in this study, the problem of over fitting is solved and the experimental results, which shown the encouraging outcomes, are achieved. The results of the experiment show that the keyword-based methods not gives better performance compare to pattern-based method. The results also indicated that removal of meaningless patterns not only reduced the cost of computation but also improves the effectiveness of the system.

## IV. PRAPOSED METHODOLOGY

The methodology which has been proposed for the solution of the problems identified in the project is as shown in the Figure 3.
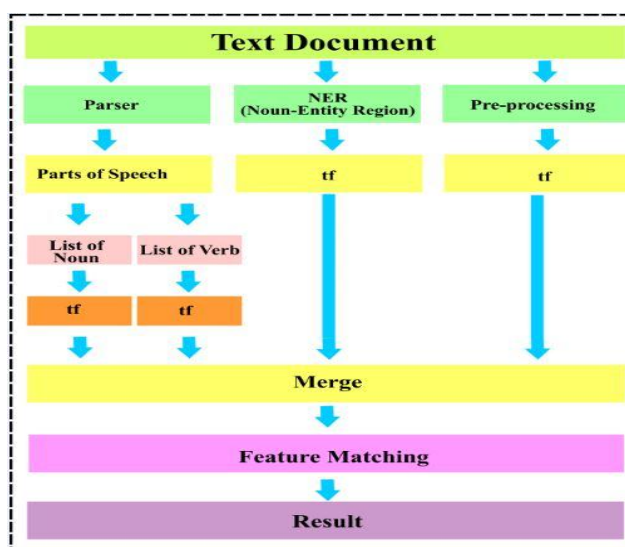


**Fig 3 Proposed Methodology**

The system is trained in three classes viz  Historical Class, Constitutional Class and Geographical Class.The system has to be trained on the basis of the documents which are related to these predefined classes. The system is trained by extracting several keywords from the document related to a particular field then the keywords that are unique with respect to each other are stored to classify the given text document.Then after the text document is uploaded into GUI to find the class to which the text document belongs.

## V. CONCLUSION

This paper's objective is to present the major techniques of location privacy. This paper surveys some of the techniques in order from the year 2012 to 2014. From the Research work, we came into a conclusion that mining the semantic information from free text document provides the enabling technology for a host to identify the class to which the text document belongs.

## REFERENCES

[1] H. Mousavi, D. Kerr, M. Iseli, and C. Zaniolo. Harvesting domain specific ontologies from text. ICSC, 2014.

[2] D. Kerr, H. Mousavi, and M. Iseli. Automatic short essay scoring using natural language processing to   extract semantic information in the form of propositions. In (CRESST Report 831). UCLA, 2013.

[3] H. Mousavi, S. Gao, and C. Zaniolo. Discovering attribute and entity synonyms for knowledge integration and semantic web search. SSW, 2013.

[4] H. Mousavi, S. Gao, and C. Zaniolo. Ibminer: A text mining tool for constructing and populating infobox databases and knowledge bases. PVLDB, 6(12):1330–1333, 2013.

[5] M. Atzori and C. Zaniolo. Swipe: searching wikipedia by example. In WWW, pages 309–312, 2012.

[6] T. Lee, Z. Wang, H. Wang, and S. won Hwang. Web scale taxonomy cleansing. PVLDB, 4(12):1295–1306, 2011.

[7] K. Raghunathan, H. Lee, S. Rangarajan, N. Chambers, M. Surdeanu, Jurafsky, and C. D. Manning. A multi-pass sieve for coreference resolution. In EMNLP, 2010