

# Architecture in Web Use Mining for Data Pre-Processing

Arti Pandey<sup>1\*</sup> Prabhat Pandey<sup>2</sup> Ajitesh S. Baghel<sup>3</sup>

<sup>1</sup>Research Scholar, Dept. of Computer Science, A.P.S. University, Rewa

<sup>2</sup>OSD, Additional Directorate Higher Education, Division, Rewa, India

<sup>3</sup>Lecturer, Dept. of Computer Science, A.P.S. University, Rewa, India

**Abstract – Data mining is a field of computer science which commingles various traditional data analysis methods with the sophisticated algorithms to process the large amount of data. There is a rapid advancement in the data collection and the data storage technologies to cumulate such kind of vast data. To extract useful information form a greater amount of data is not an easy task.**

**Keywords: Data Mining Networks, Web Mining, WUM and DSN.**

-----X-----

## 1.1 INTRODUCTION

In today life, Web has influenced the all aspect of our social life. Web has as of late turned into an effective stage for, recovering data, as well as finding learning from web data. Verifiably, the origination of finding helpful examples in Data has been given an assortment of names like Data mining, Learning Extraction, Data Disclosure, Data Harvesting, Data Archaeology, and Data Design preparing (Etzioni, 1996). According to my view the term web mining; which is worried with removing learning from web data? There has been enormous enthusiasm of Researchers towards web mining. On the premise of meaning of web mining two diverse methodologies can be proposed. One is Process Based view where the web mining is defined as a sequence of tasks and other is Data Based view which defines the web mining in terms of types of web data that was used in web mining process (Etzioni, 1996 and Cooley, et. al., 1997). In my research work we focus on data based application view which more widely accepted today.

Data based application is all the more generally acknowledged today. Web Mining is the use of Data Mining networks to concentrate learning from web data, where structure (hyperlink) or Contents/substance (genuine Data in site pages) or use Data/Usage Data (web log data) is utilized as a part of the mining procedure. On the premise of web Data three classifications of web mining are proposed, which are Web Structure Mining, Web Substance Mining and Web Use Mining. Web usage mining is the use of Data mining strategies to huge Web Data archives (Cooley, et. al., 1997). Data is

gathered in web server when client gets to the web and may be spoken to in standard configurations. The log arrangement of the document is CERN (Centre European pour la Recherche Nucleaire, or European Laboratory for Particle Physics) in Switzerland (Common log formats) (Cooley, et. al., 1999); which comprises qualities like IP address, get to date and time, ask for technique (GET or POST), URL of page got to, exchange convention, achievement return code and so on. With a specific end goal to find get to design, pre-handling is essential, since crude Data originating from the web server is deficient and just couple of fields are accessible for example revelation. Primary target of this review is to comprehend the pre-preparing of usage data. On pre-prepared Data distinctive strategies (Tsai, et. al., 2014) like factual examination, affiliation rules, sequential examples and grouping can be connected to find client get to designs.

In business, data collection is done by various ways like bar code scanners, smart cards etc. This allows retailers to collect as much as information they can find and analyze this data to enhance the business decisions. The retailers can utilize this information, along with the other business critical data. This data may include web logs from e-commerce web site and the customer service records from call centers, to help them better understand the requirements of their customers and make the business more effective and customer satisfactory.

Data mining is a field of computer science which commingles various traditional data analysis methods with the sophisticated algorithms to process

the large amount of data. There is a rapid advancement in the data collection and the data storage technologies to cumulate such kind of vast data. To extract useful information from a greater amount of data is not an easy task. The traditional data analysis tools (Tsai, et. al., 2014) and techniques are not able to serve the purpose, because of the enormous data.

The data mining techniques (Talele, et. al., 2013) can be used to support a huge variety of business applications like consumer profiling, targeted marketing, work flow management, store layout and fraud expectations. Such kind of data analysis may also provide them the names of the customers who may be beneficial for their business.

## 1.2 DATA MINING:

Data mining is all about the collecting data and then finding and extracting the useful information from such kind of large amount of data chunks. The data mining techniques are used to purge large amount of data to find the useful and interesting patterns (Cao, 2010) about the user and the business. Some of the techniques are capable of predicting the outcome of the current schemes such as whether the new customer will be willing to spend more money in the services provided by the company or not.

Data-mining is an indispensable part of knowledge-discovery. Knowledge discovery in database is the overall process to convert the raw data into some useful information. This process consists of some steps. In these steps the input data is stored in the form of a variety of formats and may reside in a centralized manner. Shown in the Figure 1.1

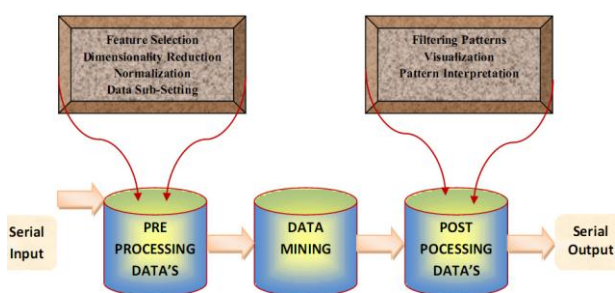


Figure 1.1 Process of Knowledge Discovery

### Figure 1.1 Process of Knowledge Discovery

Here, the data is residing in the centralized data repository to be distributed in the multiple sites. Now the data is pre-processed. The objective of data pre-processing is to transform the raw data into an appropriate format. This helps in further analysis. There are many steps involved in the steps such as fusing data from multiple sources, cleaning data to remove noise and duplicate observations, select the records and the features that are relevant to the data mining task. Since many ways are available to collect

the data and the data storage can be done in many other ways.

So, to make similarities in the data collected by these ways data pre-processing is a necessary step. After this, here come the data mining steps, which is the most important objective of this report. Data mining employs some algorithms to find out the user usage patterns, algorithms for helping business decision making and others. In the business applications, the information offered by data mining results need to be integrated with the other management tools. This integration is needed for effective marketing promotions. So this is an important goal and to achieve this data post processing step is involved. This process ensures that only the valid and the useful results are obtained by the decision support system.

### 1.2.1 CHALLENGES IN DATA MINING:

The field of data analysis is posed by some challenges (Espinosa, et. al., 2011 and Cao, 2010) These are the challenges which motivated the field of data mining. These challenges are as follows:

√ **Scalability** - In the field of data analysis, data generation and collection has grown the size of data up to terabytes or even peta bytes. Every algorithm which is handles such kind of massive data must also be capable of handling scalability. These algorithms should also involve a special data structure to access individual records in an efficient way.

√ **Heterogeneous and Complex Data** -The role of data mining in business, science, medicine and others has grown. Because of involvement of such kind of variety of data, there is a need for the techniques to handle the heterogeneous data. In the recent years there have been the emergence of the complex data e.g. collection of the web pages that contain semi-structured text and the hyperlinks etc. Techniques are developed to mine such kind of complex data objects.

√ **Data ownership and distribution** - Often it happens that the data to be analyzed, is stored in scattered locations and owned by various organized. Due to this feature there is a requirement of developing distributed data mining. But the distributed algorithms also face some hindrances. The amount of communication needed, depends on the hardware. To effectively consolidate the data mining results obtained from multiple sources is not an easy task.

√ **Non-traditional Analysis** -In the traditional statistical approach, a hypothesis is proposed, then a test is designed and then the data is collected and analyzed according to the hypotheses. But this process needs so much effort. The current scenario involves many tasks which often require generation and evaluation of many other hypotheses. The requirements of automation of hypothesis creation and evaluations have motivated some data mining techniques.

### 1.2.2 IMPACT ON OTHER DATA MINING FIELDS:

The researchers of various disciplines began to focus on developing efficient and robust tools and methods to overcome these challenges. The work done by these researchers is totally based upon the methodology, tools and algorithms that were devised earlier by the researchers.

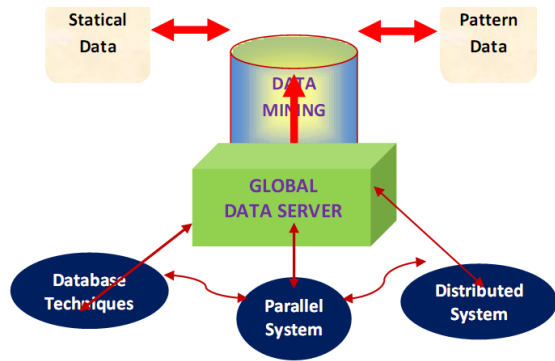


Figure 1.2 Influences the fields on data mining

Here, we explain data mining fields through the architectural diagram. The above diagram says Data Mining is by-directional connected with serial and parallel data. Another point of view data mining is internal portion of global data. Specifically data mining draws upon the ideas of sampling, estimation and hypothesis testing using statistics, search method and algorithms, modelling techniques, learning theories from artificial intelligence, pattern recognition and machine learning. Data mining is also quick in adopting ideas from other areas like information retrieval, optimization, and information theory (Kosala and Blockeel, 2000).

According to my research work, Data mining are also affected by various fields see Figure 1.2. These fields support data mining for fulfilling the purpose of collecting and extracting useful information from the large scaled databases. Data mining gets an important role in the today's computer environment. Another word it has an important research area since the amount of data available in most of the applications has become higher in size. This great amount of data should be processed to extract the

useful information and knowledge, since they are not explicit.

Apart from this Usama and et al. (Cao, 2010) defined Data Mining as "**Data mining is the process of discovering interesting knowledge from large amount of data**". All the data is residing on the www domain and all the data is being updated by many users every day. This total size of the whole documents can grow up to terabytes (Cooley, et. al., 1997). The documents on www are distributed over millions of computers that are connected with each other by telephone lines, optical fibers, radio modems and so many other technologies (Gharehchopogh and Khalifelu, 2011). World Wide Web is growing at a very higher rate in size of the traffic, the amount of the documents and the complexity level of the web sites. Due to this trend, the demand for extracting valuable information from this huge amount of data source is increasing every day (Cooley, et. al., 1999). This leads to new area called Web Mining which was first coined by Oren Etzioni (Etzioni, 1996) in 1996 in his paper, which is the implementation of the data-mining techniques to World Wide Web.

### 1.3 WEB MINING:

Web data-mining is an important sub-field of data-mining, which deals with the extraction of interesting information from the World Wide Web. Oren Etzioni in his paper claimed that web-mining uses the data-mining techniques to collect and extract the information from World Wide Web documents and services automatically. "Whether effective web mining is feasible in practice or not? , was the question which was coined by the author (Losarwar and Joshi, 2012). Web mining is the area of research, which is so much big today, due to the tremendous growth of data available on the web and the recent interest in e-commerce field. Web-mining is used to understand the customer behaviour, evaluate the effectiveness of a particular web site, and help for the success of a marketing campaign and various other decisions making for the betterment of the business (Usama, et. al., 1996 and Sharma, et. al., 2011). There is one simple definition of web mining; which is "Web Data Mining is the application of data mining techniques, to extract interesting and potentially useful knowledge from web data. It is normally expected that either the hyperlink structure of the web or the web log data or both have been used in the mining process". (Usama, et. al., 1996).

#### 1.3.1 EVOLUTION OF WEB MINING:

The information gathered through Web mining is evaluated (sometimes with the aid of software graphing applications) by using traditional data mining parameters such as clustering and classification, association, and examination of

sequential patterns. The evolution and emergence of web mining is as follows:

- √ **Data Collection (1960)** is portrayed by its review nature and static Data conveyance.
- √ **Data Access (1980)** is portrayed by its review nature and dynamic Data conveyance at record level.
- √ **Learning Discovery Database (1989)** is the non-paltry extraction of understood, beforehand obscure and possibly valuable Data from databases.
- √ **Data Warehousing and Decision Support Networks (1990)** is described by its review nature and dynamic Data conveyance at different levels.
- √ **Data Mining (2000)** is portrayed by its forthcoming nature and proactive Data conveyance at any level.
- √ **Web Mining (2010)** is described by its planned nature and hyperactive Data conveyance at any level. Web mining in perspective of Data mining have three basic operations:
  - √ **Clustering:** discovering characteristic gathering of clients, pages and so on.
  - √ **Sequential get to Pattern Analysis;** the approach of back to back blue print exposure endeavours to decide between sessions blue print so that the comportment of a gathering of tokens is sought after by another token in a period administered gathering of sessions or portions. By using this technique, organize advertisers can expect next force blue prints which will be helpful in laying adverts coordinated at specific customer gatherings. Different sorts of mainstream investigation that should be possible in back to back blue prints comprise obviously examination; adjust point spotting, or comparability investigation.
  - √ **Associations Rule:** which URLs have a tendency to be asked for together. Affiliation govern era can be utilized to relate pages that are regularly referenced together in a solitary server session.

### 1.3.2 ASSERT IN WEB MINING:

- √ **To find relevant information:** The users involve internet, in their searches and for gaining knowledge. However today's searching methods have problems like low precision which is due to the irrelevance of

many of the searching results. These problems result in the difficulty to find the relevant information. Another problem is low recall which is due to inability to index all the information available on the web (Espinosa, et. al., 2011).

- √ **Getting new information from the data available on the web:** This problem is fundamentally a sub-problem of the first one. Above problem is query triggered process but this problem is a process which is triggered data. In this problem, the useful information is collected from the data stored (Gharehchopogh and Khalifelu, 2011).
- √ **To learn consumer's needs:** This problem is all about that what a user needs and searches often. Inside this problem there are sub problem such as customization of the information to the intended group of consumers. This problem is more related to the web site design management and marketing (Sharma, et. al., 2011).

### 1.3.3 ASSESSMENT OF WEB MINING: (Kosala and Blockeel, 2000)

The various tasks of Web Mining are as follows:

- √ **Asset Finding:** This is the task of retrieving the desired records proposed for web documents. The word resource finding means, the process of retrieving the data which is either online or offline, from the sources available on the web, such as electronic newsletters, the text contents of HTML documents obtained by removing HTML tags and also the manual selection of web resources (Losarwar and Joshi, 2012).
- √ **Data Selection and Pre-handling:** Automatically choosing and pre-preparing particular Data from assets recovered from the web sources is the heartbeat of this phase. It is like a transformation process of the original data retrieved in the IR process into a useful data-set. These transformations can be pre-prepared such as stop words, stemming or a pre-processing aimed at obtaining the desired representation such as finding phrases in the training corpus, transforming the representation to relational or first order logic form (Sharma, et. al., 2011).
- √ **Speculation:** Discover general examples at individual site destinations and in addition over different locales site. Machine learning or data mining techniques are used in the

process of speculation (Kosala and Blockeel, 2000).

✓ **Investigation:** Validation as well as translation of mined examples.

**1.3.4 WEB MINING TAXONOMY**

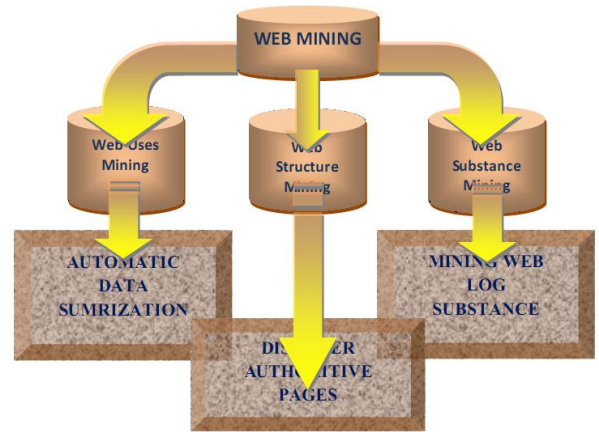
The term Data mining is characterized as the programmed extraction of unidentified, helpful and justifiable examples from expansive databases. So as to build the execution of Website, the fundamental thing is great web architecture (Usama, et. al., 1996). The interests of the clients help in planning better Websites. Web mining is utilized to recover, extricate and assess Data for Data disclosure from archives on Web. Web mining comprises of three sub fields: Web Content or Substance mining, Web structure mining and Web Usage or Use mining (Cooley, et. al., 1999).

- ✓ Web content Mining manages the pattern of Data which is valuable from the Web Data or Reports.
- ✓ Web Structure Mining mines the hyperlinks structure inside the web itself. The Structure speaks to the chart of the connection in a Website.
- ✓ Web Usage Mining mines: the data at log document put away in the web server.

I'll present these fields through follows diagram. See Figure 1.3, it's having their own importance according to the need of implementation. The Web mining is filtered into three different ways; Web Uses Mining, Web Structure mining, and Web Substance Mining continuously. As we know that WUM is the application of data mining techniques (Srivastava, et. al., 2000). Is to discover interesting usage patterns from Web data in order to understand and better serve. Web structure is to discover useful data from hyperlink. The web structured is analyse link structure of the web, and also identify more preferable documents for end users

The Structure mining or structured data mining is the process of finding and extracting useful information from semi-structured data sets. And Web substance mining is extracting something useful or valuable from a baser substance, for example; mining gold from the earth. And, also is used to understand customer behaviour, evaluate the effectiveness of a particular Web site, and help quantify the success of a marketing campaign. Automatic data summarization is part of machine learning and data mining. The main idea of summarization is to find a subset of data which contains the information of the entire set. Such techniques are widely used in industry today. For example Search engines; others

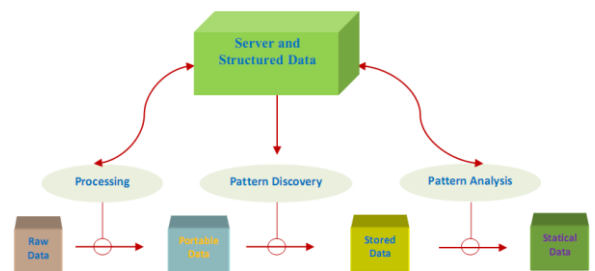
include summarization of documents, image collections and videos.



**Figure 1.3 Advanced Block Architecture of Web Mining**

**1.4 WEB USAGE MINING (WUM):**

The term web usage mining was presented by Robert Cooley et al. in 1997 according to which “Web Usage Mining is the programmed disclosure of client get to designs from web servers”. The procedure of revelation and examination of examples concentrates on client get to Data (web use data). Web perusing conduct of clients is caught by Web use Data from site. In our unique situation, the usage data is Get to sign on server side that keeps data about client route. According to my view global server stored structured data for long time I its database. Which can be analysed and proceed by the end users. They use that stored data as Raw material, processing this data in form of innovative data. In another way we say that end users can use predefined pattern for creating new ideas or new data. See this approach in Figure 1.4; this diagram demonstrates the distinctive periods of web use mining.



**Figure 1.4 WUM Process Block Architecture**

At last, according to my concept, there is analysed information on web pages visits that are saved in log files of Internet servers in order to discover the previously unknown and potentially interesting useful patterns. Web usage mining is described as applying data mining techniques on Web access logs to optimize web site for end user interest i.e. Web Uses Mining.

### 1.4.1 DATA SOURCES FOR WUM:

The Data in log is gathered from sources like server side, customer side and intermediary servers et cetera. The data accumulation step incorporates different data sources. The Primary wellspring of data in web use mining is the log at server. There are some extra data source are likewise use for some client and some application which incorporates sign on customer side and Proxy side log (Cooley, et. al., 1999). In Log at customer side, usage data can be followed likewise on the customer side. In many regards, gathering route data at the intermediary level and at server level is same. The fundamental distinction is just that intermediary servers gather data of client gatherings getting to huge gatherings of web servers. Data which is utilized for web use mining can be gathered at three unique levels (Ou, 2011).

- √ **Server Level:** The server stores Data with respect to ask for performed by the customer. Data can be gathered from numerous clients on single site. A Web server log explicitly records the browsing behaviour of site visitor, thus an important source of WUM. The multiple users access data recorded in server logs reflects the (possibly concurrent) of websites. These log files can be stored in various formats such as Common log or extended log formats. The other kinds of usage information such as cookies and query data in separate logs can also store in the Web server.
- √ **Proxy Level:** A Web proxy acts as an intermediate level of caching between client browsers and Web servers. Proxy caching can be used to reduce the loading time of a Web page experienced by users as well as the network traffic load at the server and client sides. The performance of proxy caches depends on their ability to predict future page requests correctly. Proxy traces may reveal the actual HTTP requests from multiple clients to multiple Web servers. This may serve as a data source for characterizing the browsing behaviour of a group of anonymous users sharing a common proxy server.
- √ **Client/Browser Level:** - A Client side data collection can be implemented by using a remote agent (such as Java scripts or Java applets) or by modifying the source code of an existing browser (such as Mosaic or Mozilla) to enhance its data collection capabilities. The implementation of client-side data collection methods requires user cooperation, either in enabling the functionality of the Java scripts and Java

applets, or to voluntarily use the modified browser.

These methods will collect only single-user, single-site browsing behaviour. A modified browser is much more versatile and will allow data collection about a single user over multiple Web sites. The most difficult part of using this method is convincing the users to use the browser for their daily browsing activities. This can be done by offering incentives to users who are willing to use the browser, similar to the incentive programs offered by companies such as NetZero (Cao, 2010) and All Advantage (Cooley, et. al., 1997). that reward users for clicking on banner advertisements while surfing the Web.

### 1.4.2 STAGES OF WUM:

Web usage mining uses data mining methods and techniques on extensive web log archives to find learning which is helpful about behavioural example of client and further more site usage insights that can be utilized for different web composition assignments [11, 12, and (Srivastava, et. al., 2000). The four stages under web use mining are:

#### STAGE I-

- I. **Data Pre Processing** In this stage log files are transformed into a form that is suitable for mining. This is done on crude Data which show in log document wrapping up of Data cleaning, client recognizable proof and session distinguishing proof. The data accessible in the web is Varied and unstructured. Consequently, the pre-preparing stage is a required for finding designs. The reason for this is to change the crude data into a gathering of client profiles. Data pre-handling is essential and this prompted different calculations and heuristic strategies for it. The following are preprocessing tasks that have been identified.
- II. **Data Cleaning** Data Cleaning is a procedure of expelling things which are superfluous, for example, jpeg, gif documents or sound records. The enhanced data quality likewise enhances the examination on it. In the event that a client demands to see a specific page alongside server log sections the scripts and illustrations are downloaded with a HTML record. Additionally Check the Status codes in log sections for effective codes.
- III. **Client Identification** The recognizable proof of individual clients who get to a site is an imperative stride in web usage mining Process. Different techniques are to be taken after for this. The easiest technique is

to a lot particular client id to unmistakable IP addresses. In the event that the client's IP address is same as past passage and client specialist is distinctive, then the client is accepted as another client. On the off chance that the page that is asked for is not specifically reachable from any of the pages till gone to by the client (Tsai, et. al., 2014), then the client is recognized as another client in a similar address.

- IV. **Session Identification:** The arrangement of pages gone by a similar client inside the term of one particular visit to a site is considered as a session of client. There are more than one session related with same client moreover. The one technique relies on upon time and another on route in web topology utilized for recognizable proof of sessions. In Time Oriented Heuristic (Etzioni, 1996), there are two strategies in which one technique in view of aggregate session time and the other in light of single page stay time. The arrangement of pages which are gone to by a client at a particular time is called page seeking time. The second strategy relies on upon stay time on page which is ascertained with the contrast between two timestamps. These strategies are not solid since clients may include in some other work in the wake of opening the website page.
- V. While in Navigation-Oriented Heuristic, the thing which is considered is website page availability. In the event that a site page is not associated with page which is opened beforehand in a session, then it is considered as another session. Both the techniques are utilized by numerous applications.
- VI. **Path Completion:** Some accesses are not captured by the web log file and this may result in *incomplete paths* (requests made to a page not directly linked to the last requested page) in the log data. This is probably due to the use of local caches or proxy servers. Some solutions include *cache busting* (using cache specific headers to stop page caching) or the use of navigation history (Pitkow, 1997).
- VII. Path Completion is the procedure to finish the get to way of client utilizing URL and referrer page get to way entire structure of pages and connections of pages that are gotten to by a client. Charts are utilized to speak to the ways of client get to. Charts are utilized to speak to the way finish handle. Each hub is utilized to speak to a page or a

site and edges speak to the connections between the pages or sites.

- VIII. **Transaction Identification:** Once a session of completed paths is determined, page references must be grouped into logical units representing Web transactions before any mining can be carried out.

## STAGE II- APPLYING MINING FOR PATTERN DISCOVERY:

There are four basic mining techniques that can be applied to Web access logs to extract knowledge, but we will focus on algorithms based on association rule mining (ARM) and sequential Pattern Mining (SPM) because of their complexity and applicability. These are the few techniques;

- √ Sequential-pattern-mining-based: recognize the discovery of temporally ordered Web access patterns.
- √ Association-rule-mining-based: determine the correlations among Web pages.
- √ Clustering-based: Grouped the users with similar features.
- √ Classification-based: Grouped the users into predefined classes based on their features.

The above techniques help in the design of better Websites as well as in the development of effective marketing strategies.

## STAGE III-

The third stage of web usage mining Process is Pattern Analysis or example analysis. The examples are found in this stage perform likewise the measurable investigation, affiliation rules, bunching, design coordinating et cetera. The examples which are dug are not reasonable for elucidations. So it is critical to deal with examples or guidelines which are not intriguing from the set found in the example revelation stage. The devices are given to help the change of data into learning in this stage. The correct investigation is administered by the application for which web mining is finished. The SQL is the most well-known strategy for example investigation. While another strategy is to load usage data into a data shape so as to perform OLAP operations. Once pattern were found from web logs, the guidelines or patterns which are not fascinating are shifted through.

**STAGE IV-****APPLYING MINING RESULTS:**

The last stage of WUM involves the detection and transformation of mining results into useful actionable tasks such as:

- √ Re-design Websites so that correlated pages are found together.
- √ Improve access time by pre fetching pages frequently accessed sequentially.
- √ Improve caching by storing pages frequently revisited.
- √ Enhance surfing experience by relocating pages in such a way that users need not visit unnecessary pages to get to their desired pages.

All the four stages are appeared through the accompanying. When exchanges of client have been distinguished, methods of data digging are performed for example disclosure in web use mining process (Ou, 2011). These strategies speak to the ways that regularly show up in the data mining study, for example, disclosure of affiliation guidelines and successive examples and bunching and arrangement and so on. Characterization is an administered learning process. In this, the data mapped into one of a few predefined classes, it should be possible by utilizing inductive learning calculations, for example, guileless Bayesian classifiers, choice tree classifiers, Support Vector Machines and so forth. Bunching is a strategy of collection clients which display comparative perusing designs. Such examples are valuable for deriving client tally keeping in mind the end goal to perform showcase consider in E-Commerce or E-Learn give customized web substance to pages. By utilizing this strategy, web advertisers can anticipate future visit designs which can help in setting ads gone for certain client gatherings.

**1.4.3 APPLICATION OF WEB USAGE MINING:**

It is the application which utilizes diverse data mining networks to break down and separate intriguing examples of client's use and interests over data on web. The usage data comprises of client's conduct while perusing on web (Huan, et. al., 2010). This movement includes finding the examples naturally from at least one network has. Frameworks that use this application render and accumulate tremendous greater part of data, ordinarily rendered mechanically by network has and kept up in host logs. Frameworks inspect this data that serves to discover the client's advantages, cross showcasing plans and limited time crusading techniques and so on.

Web mining is a clever investigation of Web data (Huan, et. al., 2010). WUM is the procedure which extricates "intrigued" blue prints from the network data. The network data comprises of network host get to log, entryway have log, program log, customer enrolment data, and customer's session. In this we primarily utilize web log as data source. So we utilize the idea of web log mining rather than WUM. The procedure of web log mining is as per the following:

- I. **Data Readiness;** is the primary phase of web log mining. The crude data is changed over into the data with which design disclosure could bargain. It incorporates data cleaning, client acknowledgment, session acknowledgment, way supplement, exchange acknowledgment et cetera. Web log data pre-handling directly affects the accuracy or models and example rules which are found in the following stage.
- II. **Designs revelation;** utilizing different strategies, we endeavour to discover models and example standards of client's get to conduct. Normal innovations are consecutive examples, affiliation tenets, bunching, and order et cetera.
- III. **Design investigation;** it is utilized to concentrate significant intriguing examples from every single existing model. In a large portion of the cases, web use mining can discover every one of the models and principles.
- IV. **Pre-handling;** it incorporates the strategies like cleaning the data, customer acknowledgment and stage acknowledgment. These techniques are utilized to the real web log documents to gain finish web get to sessions. Data cleaning is considered as site particular process which incorporates vital assignments like joining the logs from a few servers and making lumps of the logs into data things. Be that as it may, the illustrations document solicitations are expelled from the log records after pre-preparing.
- V. The principle aim of the pre-processing procedure is to pre-handle the strict network logs to discover complete network get to sessions. While using the network have logs, every customer' s get to undertakings and works completed by the customer of a site are commented around the network host of the site. Every customer get to data incorporates the customer web convention address, appeal to time, required Uniform Resource Locator, Hyper Text Transfer Protocol status code, and so on. Customers



are considered overall gathering as the web convention locations are not identified with every customer's conspicuous deceivability data. Ordinarily, organize logs might be viewed as a gathering of continuous of get to tokens from one client or stage in a day and age expanding request.

**VI. Data cleaning;** Data cleaning is a site particular stride that includes unremarkable errands, for example, combining logs from numerous servers and parsing the sign into data fields. Regularly illustrations document solicitations are stripped out at this stage. This is effectively done by checking for record names postfixes, for example, " GIF " or "JPG" Graphics documents can be left in the data all index and moved up into site visits in a later pre-handling venture with no loss of all inclusive statement or demand for some other record which might be conceded into a network page: or notwithstanding marine session performed by robots and network creepy crawlies. While request of for graphical subjects and documents are tender to destroy robot and network creepy crawlies nautical blue prints must be blue print must be expressly. This is typically practiced for instance by referring to the inaccessible host, by referring to the specialist, or by guaranteeing in the entrance to the rpbpts .txt document. HTTP status figures are used to speak to the win or lose of the called for issue. The records with figures among 200 and 299 are viewed as profitable records, and remaining are expelled from the networks logs.

**VII. Customer and session acknowledgment;** For analyzing customer get to direct, unparalleled customers must be perceived. As specified some time recently, customers are viewed as mysterious in most network hosts. We can modify the client acknowledgment technique to customer Internet Protocol acknowledgment. In an alternate dialog, petitions from a similar Internet Protocol address can be considered as from a similar customer and kept in a similar bunch under that customer. To perceive customers all the more accurately, some other data from the network logs might be helpful. The agentive part enlisted in network logs catches data on the client program on Formal Based Concept Analysis. At that point usage of our recommended WUL-unearting calculation to uncover is the most potential and utilitarian gathering of connection get to blue prints from the Network Usage Lattice. The vantage of the recommended WUL-grounded technique is that it can create substantially less number of

connection get to blue print standards without trading off much on lineament for network individualization applications when compared with the Apriori-Grounded algorithms [36]. Web Usage Mining (WUM) techniques have been connected to numerous constant down to earth applications [30] including the followings:

**VIII. Personalization;** Web usage unearthing methodologies can be used to supply individualized network program involvement. For instance it is conceivable to anticipate the program lead in exacting time by likening the present marine blue print with particular blue prints that were evoked from past network log. In this field, recommendation frameworks require the most common application; its principle objective proposes concerning connections to items that could worry to a number of the customers. Individualized Site Maps are a case of proposal framework for connections recommended an adaptive technique modify the item index associating to the evaluated customer perceivability. A technique to join handle cosmologies into the individualization method grounded on web use uncovering is recommended in conceding a calculation to assemble data base level mass visibilities from a gathering.

**IX. Network Improvement;** Rate of value and another quantifiable quality property are extremely fundamental to customer satisfaction from administrations like databases, networks and numerous progressively and same sort of the qualities are anticipated from the customers of web administrations. Web usage mining gives the thought to comprehend the web stack conduct that might be valuable to build rules for web reserving network transmission, stack adjusting or data conveyance. The real concern may the arrangement for security for electronic administrations especially as e - trade keeps on developing at an exponential rate. Web use mining is additionally valuable examples that are useful to distinguish interruption, wrongdoings, endeavoured break-ins, and so forth., Some models are recommended to anticipate reality, to the two worldly too spatial, in the website pages which are required from a particular customer or a group of customers who access from the comparable intermediary server. The parameters are additionally subject to the area of the server which are required to choose consummating and getting approaches for the intermediary server. When utilizing the greater amount of the steadily changing substance has

diminished the benefits of putting away at the customer level and server level.

**X. Site Modification;** The quantifiable nature of the a site which incorporates the pulling in power, with regards to both substance and skeleton is extremely fundamental to the greater part of the applications, for example, an index for the items inventory fore-trade. WUM renders expounded resubmit on customer direct, providing the network site planner with data on which to ground re-project conclusions. Arrange use data supplies an opportunity to end up noticeably each site into a progressing useable trail. As their data is not as refined as the data that can be gathered from a formal useable investigation with recordings and edified percipients, Network use data are cheap and abundant. Get to times and measures of lost can be computed naturally rather than physically. At the point when the results of any of the recommendations could prompt re arranging the structure and subject of a site focuses on mechanically changing the structure of a site grounded on use blue prints revealed have logs. Constellating varlets is connected to discover which varlets ought to in a flash relate.

**XI. E-Business Intelligence:** Data on how clients are utilizing a site is basic for advertisers of web based business organizations. Buchner has introduced an data procedure keeping in mind the end goal to find showcasing insight from web data. They characterize a web log data hypercube that combines web usage data alongside advertising data for web based business applications. Four particular strides are recognized in client relationship life cycle that can bolster by the data revelation methods: client fascination, client, cross sales and client flight. Web server logs use to produce convictions about the get to example of website pages at given site .Algorithms for finding intriguing guidelines in view of the surprise are additionally created.

## 1.5 CONCLUSSION:

In this research paper I've study the web usage mining process and its architecture for analysis. It is expected that the work will be contribute in suggesting some improvements algorithms as well as the implementation of Apriori Algorithm for data Pre-processing in web usage mining. The effective algorithm will be proposed with the improvements. The proportional analysis on various association rule mining using on the clustering and its applications in the web log files.

## REFERENCES:

- B. Roberto Espinosa, J. Zubcoff and J. Maz'on (2011). "A Set of Experiments to Consider Data Quality Criteria in Classification Techniques for Data Mining", *International Conference on Computational Science and It's Applications" (ICCSA) 2011, Part II, LNCS 6783*, pp. 680–694.
- C. Tsai, C. Lai and M. Chiang (2014). "Data mining for internet of things: A survey," *IEEE Communication Surveys & Tutorials*, Vol. 16, No. 1.
- D. V. Talele, V. Dipali. and C. D. Badgujar (2013). "A Literature review of Opinoin Mining From Online Customers' Feedback and It's Applications Domains", *Asian Journal of Computer Science & Information Technology*, Vol.11, pp. 301-305.
- Definition of Data Mining", Available: <http://www.britannica.com/EBchecked/topic/1056150/data-mining>, May 2012.
- F. S. Gharehchopogh and Z. A. Khalifelu (2011). "Analysis and evaluation of unstructured data: text mining versus natural language processing" *Application of Information and Communication Technologies (AICT), 2011 5th International Conference on*. IEEE, pp. 1-4.
- F. Usama, P. Shapiro, and S. Padhraic (1996). "From Data Mining to Knowledge Discovery in Databases", *American Association for Artificial Intelligence*, USA, pp. 12-17.
- K. Sharma, G. Shrivastava, and V. Kumar (2011). "Web mining: Today and tomorrow", *Electronics Computer Technology (ICECT), 3<sup>rd</sup> International Conference*, Vol. 1 IEEE.
- L. Cao (2010). "Domain-driven data mining: Challenges and prospects, " *Knowledge and Data Engineering", IEEE Transactions*, Vol. 22, No. 6, pp. 755-769.
- M. L. Huan, R. Setiono and Z. Zhao (2010). "Feature Selection: An Ever Evolving Frontier in Data Mining", *Workshop and Conference Proceedings 10*, pp. 4-13.
- O. Etzioni (1996). "The World Wide Web: Quagmire or Gold Mine". *Communications of the ACM*, 39 No.11, pp. 65-68.
- R. Kosala and H. Blockeel (2000). "Web Mining Research: A Survey", *ACMSIGKDD*, Vol. 2, No. 1, pp. 1-15.

Robert Cooley, Bam Shad Mobasher, and Jaideep Srivastava (1997). "Grouping Web page references into transactions for mining World Wide Web browsing patterns", Knowledge and Data Engineering Workshop, New port Beach, CA.IEEE, pp. 2-9.

Robert Cooley, Bamshad Mobasher, and Jaideep Srivastava (1999). "Data preparation for mining World Wide Web browsing patterns", Journal of Knowledge and Data Network, pp. 1-27.

Srivastava, J., Cooley, R., Deshpande, M. and Tan, P.N. (2000). "Web usage mining: discovery and applications of usage patterns from Web data" ACM SIGKDD Explorations Newsletter, Vol. 1, No. 2, pp. 12-23.

V. Losarwar and M. Joshi (2012). "Data Preprocessing in Web Usage Mining", *International Conference on Artificial Intelligence and Embedded Systems (ICAIES'2012)*, pp. 15-16.

Web mining definition available:  
[http://en.wikipedia.org/wiki/Web\\_mining](http://en.wikipedia.org/wiki/Web_mining).

Z. Ou (2011). "Data Structure and Effective Retrieval in the Mining of Web Sequential Characteristic", *International Conference on Electronic & Mechanical Engineering and Information Technology*, pp. 3551-3554.

---

### **Corresponding Author**

**Arti Pandey\***

Research Scholar, Dept. of Computer Science,  
A.P.S. University, Rewa

**E-Mail – [India.aarti.tiwari10@gmail.com](mailto:India.aarti.tiwari10@gmail.com)**