# Detailed Survey on Recommendation System with Classification Technique in Web Usage Mining

Neha Khandewal, R.L Yadav

Research Scholar, Associate .Proffesor

Dept Computer Science and Engg

Kautilya Institute of Technology & Engineering ,Jaipur

*Abstract*— As all data is accessible on the Internet however it is difficult for each client to discover applicable data in a short span of time. So as to beat this issue proposal framework presented in the Web world. In this paper, the issue and different methods are clarified for suggestion models. The serious problem of various online web sites is the overview of many decisions to the altered customers at the same time. This normally results in a tedious assignment in determining the accurate object or information on the site. The client's extant conspiracy trusts on the directional conduct It cooperates with clients directly in their interventions and encourages them to take some time in limited quantities in a limited capacity. Further examples obtained through data collection strategies did not work well in the future of future perusing designing due to the low coordinated rate of the guidelines and the supervision of the client. Web Usage Mining (WUM) finds challenging examples of use from web-based information and promotes the need for web-based applications in a major way.

*Keywords*— *Data Mining, Web Usage Mining, Recommendation System, Really Simple Syndication, Classification, K-NN Algorithm.*

## I. INTRODUCTION

Data Mining is a procedure of examining data from specific summarizing and perspectives the final outcome as helpful information. It has been describing as "the non-trivial procedure of new, potentially helpful and ultimately conceivable patterns & identifying valid in data". The objective of a data mining exertion is typically either to make an illustrative model or a prescient model [1].

Web mining is the utilize data mining execution towards without human intercession find and mine data from Web reports and administrations. The Web mining examination is a meet up with research zone from a few research networks, for example, database, IR, and AI look into network particularly from machine learning and NLP [2].

Web usage mining(WUM) is the utilization of data mining procedure to consequently find and concentrate helpful data from a web site. Lately, there has been a hazardous development into the quantity of examines into the field of web mining, explicitly of WUM. The client access and route example model are removed from the verifiable access information recorded in the client's RSS address URL document, utilizing proper data mining strategies. The K-Nearest Neighbor arrangement strategy was utilized online and in Real-Time to endeavor WUM system to distinguish customer's/guests clickstream information coordinating it to a specific client gathering and suggest a custom-made perusing choice that addresses the issue of the particular client at a given time [3].
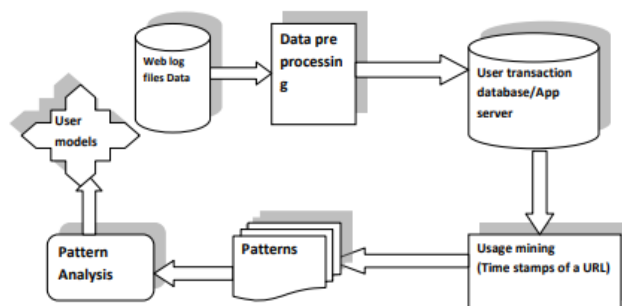


Fig. 1. Web usage mining process

## II.  WEB USAGE MINING

Web Usage Mining is a procedure by which extracting useful data after server logs, for example, Web usage mining is the procedure of definition out what clients are searching proceeding the Internet. Some users may only show interest in multimedia data, sometimes with interest in multimedia data. [4] Web Usage Mining is the use of information-based applications to find exciting usage designs after web information. One or more Web sites or web resources can be referenced or collected consequently of client communications or auto detection, analysis of the related clues & web application mining. Capture, model & consider the behavior & profiles of clients engaged in a Web site. The patterns found in the collection of pages, materials, or resource materials that are often accessed by groups of users with common needs and interests are usually found. Data uses data or identity or data using their browsing behavior proceeding the Web site.

Caching and proxy servers are typically a problem for individual users, service sessions, episodes, etc. Detail and association of certain pre-processing approaches aimed at web usage information.

1) **Data Collection:** In this step, the data will be collected from web servers or clients visiting the web site.

2) **Data Preprocessing:** This is especially important for session identification, data cleaning, & user identification.

3) **Pattern Discovery:** In this step, information is distinguished based on information such as mechanical techniques, cluster, classification, association roulette discovery, etc.

4) **Knowledge Post-Processing:** In this final phase, understanding is understood into a procedure that is comprehensible to humans, e.g. With reports, or visualization methods. In addition to these techniques, the results posted by the post can also be included in a web-personalization section.

## III. RECOMMENDATION SYSTEM

The recommended systems are beneficial for both service providers and users, and the information distribution system that deals with system issues. The transactions for finding and choosing items in an online shopping environment will be reduced. Systems of decision-making and quality improvement recommendations are also proven. In e-commerce settings, recommended companies are raising revenue, which is an effective way to sell more products. Supporting systems in science libraries allow users to navigate through catalog searches.
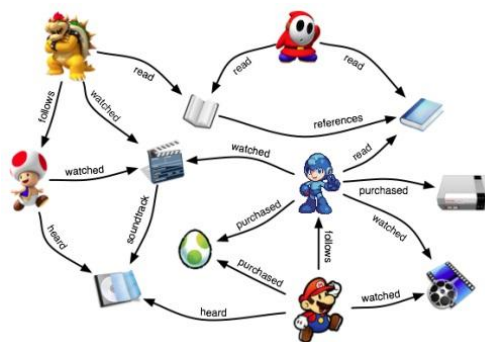


Fig. 2. Recommendation system

A.     **Phases of Recommendation Process**

- *Information collection phase:* It gathers the appropriate data by clients to create a client profile or model aimed at predictive challenges comprising user access, behavior or contents of properties accessible by the client.

- *Learning phase:* A Study Algorithm applies to filter and exploiting client's aspects after feedback collected into a data repository step.

- *Prediction/recommendation phase:* It proposes what type of items the user has selected, or predicts. It is directly based on data collected in the data collection phase based on memory-based, model-based or user-run processes.

### B. Classification of Recommendation System Approaches

Recommendations are generally categorized based on their rating approach [5]:

- **Collaborative Filtering System:** Collaborative filtering (CF) systems act as a number of users by collecting user feedback in rating formats of the given domain and creating rating behavioral similarities.

- **Content-based System:** The content-based recommendation mentions that the client is based on the item description and a user's interest profile. Web pages, TV such systems are used to refer to events, news articles, etc.

- **Hybrid System:** The hybrid recommendations are the techniques that combine multiple recommendations to achieve a synergy between.

### C. Problems in Recommendation Systems

Various technologies used on a recommended system experience some disruptions based on basic issues [6].

**Sparsity Problem:** In any recommended structure, the number of ratings previously received is so small associated with the number of ratings to be forecast. The advantage of ratings after a small number of instances is significant.

**Cold-Start Problem:** Novel objects & new clients are making a big task for the recommended structures. These issues are called interactive cold-start issue. The 1$^{st}$ of these concerns is into the premix of Collective Filtering. Cannot recommend a single item if some users did not rate this earlier.

**Scalability:** As the number of users and IT companies increases, additional resources are needed to improve the systems and to set up recommendations. Most of the resources use the same information materials as well as to determine users with similar tastes.

## IV. REALLY SIMPLE SYNDICATION

Really Simple Syndication (RSS) aggregator is XML-based layouts for syndication of high listing links, helping users to decide their links by regularly updating digital content. It's officially used to gather news headlines. There are two types of RSS-aggregator, including browsers built RSS devices & user-based RSS devices. The RSS "Really Simple Syndication or Rich Site Introduction" is a method to get updates deprived of visiting separately. Internet clients may utilize an RSS feed reader or aggregator to accomplish their particular RSS feed contributions. Simon claims that RSS feeds cannot be a system to display headings and show some stories. Users control how long it takes to users' sites to see how many people are viewed at a time [7]. Password is the major restriction of Web sites - the site should access the site to perceive what's original. The RSS, such as Ajax, Flash & Silver Lite uses, brings this kind of popularity to this particular purpose, providing an active way to attract traffic and personal consideration. [8].

## V. CLASSIFICATION

Each day, the automatic classification of documents in previous categories attracted several researchers. The supervised, unprivileged & semi-supervised approaches are utilized to redesign the documents. In the past decade, we have seen the baseline classification, decision tree, neighboring neighbor, support vector machines, neural networks, and Rocchio's.

### A. K-Nearest Neighbour

The nearest neighbors are a virtual supervised machine learning algo that stores completely accessible cases. Novel cases built on a similar measurement (such as distance functions) may be grouped. K-NN points in close proximity (documents) are from a similar class. The algo implies all the training examples by analyzing a fixed number of neighboring neighbors (K) by using some similarities of university distance.

### B. Support Vector Machine

Initially, support for creating an optical binary (2-class) classifier was developed by vector machines (SVM), but later retracting & clustering issues were developed. The operating belief of SVM is to discover a high-plane (linear/linear) that exploits margin. SVM is a fractional case for kernel-based approaches. It connects aspect frames to a high discrete distance with a kernel function that builds a stunning hyper-platform that matches the optimal linear discriminating function or training data on this space. The kernel is not well-defined into the SVM case. In its place, we need to define the distance among any two opinions into the hyperspace.

### C. Naïve Bayes

Simplified Bees Classics are a balancing act based on the theoretical principle of freedom and helplessness. This is one of the most simple text organization methods using several uses on email spam detection, particular email sorting, and classification of documents, sexually explicit content detection, sentiment detection & language detection sentiment detection. Trials demonstrate that this algo will perform well in numerical & text information. As it has very low computational intensity (CPU and memory), it has done better in other technologies such as simple bumps, wood trees, max entropy and support vector machines. Training information. The speculation of systematic individuality is related to actual world information. It has good performance and its performance.

### D. Neural Networks

Remote networks may be utilized to create difficult connections among I/Ps & O/Ps to discover designs into the information. When mining networks are used as a tool, data mining companies collect data from datasets. An NRC network is a chain of unit units, while input units are usually represented by words, and o/p unit (s) denotes this classification. To classify a text document, its word waiters are given to i/p units; the activities of these units are put forward by the n/w, & the o/p unit (s) are increasing to a certain extent and the determination of the sorting decision.

### E. Rocchio's Algorithm

Rocchio's algorithm was built on the concept feedback concept method available in the 1970s, as it was popular among smart data reclamation structures after the Smart Data Reclamation Structure. This algo has a prototype vector based on every class. The standard vector of a prototype vector in all training line vectors of Class $c_i$. The text document the text between each prototype capabilities and assigns the text document to the maximum class of classes. The algorithm is based on a general idea to indicate which documents are relevant or irrelevant for most users. This algorithm is fast-tracking & easy to implement. Though simple to execute, this algo comes after bad organization precision. Performance in the column of permanent alpha and beta ranks an important part of its performance.

## VI. K-NEAREST NEIGHBOR

As the name algorithm suggests, a novel tube classification is associated with the nearest tube. Starting the classification process on KNN begins by the data set. The information design of some of the special attributes that describe the data set. The data set is distributed into 2 groups: the training set & the test set. The training set algorithm will also be provided when the test set is utilized to detect algorithm precision. In KNN algo, [10], Trouble Tapes may be observed such a group of data points into an-dimension, i.e. the division of the information set may be complete with approaches, for example, hot-out technique, random sampling, cross-validation. In space, n sets are the sets of n attributes that explain the data set. While an unknown toll aimed at organization arrives, it will have to discover the nearest information points into the n-dimensional area. For instance, K. To find the longest distance to know the tube for Euclidean distance, Minkowski distances, and Manhattan distance, The Euclidean distance among 2 information tuple X & Y is specified below: KNN is the most commonly used algorithm when using different algorithms for classification. The KNN algorithm is very easy to implement and give good results. KNN algorithm does not require prior knowledge of data for classification. It stores all training tutorials provided for input without inputting anything. All calculations are done during a test tuple classification. For instance, K. To find the longest distance to know the tuple for Euclidean distance, Minkowski distances, and Manhattan distance, The Euclidean distance among 2 information tuples, X & Y is specified further down:

## VII. LITERATURE REVIEW

Rafael de Oliveira Werneck et al. [2018] the Kuaa suggests a workflow-based framework aimed at enterprise, distribution & implementation of an automated machine-learning learning experiment. This framework is designed to design a framework for analyzing machine learning solutions because the aspect descriptors, normalizers, classifiers, & fusion techniques help into a broad variety of activities related to machine learning. Users who use Kuaa with the help of machine-learning workflows. Utilize of Recommendations agrees clients recognize, estimate, & revise successful solutions that were previously defined. In the implementation of the recommended services, we suggest the utilize of similar measurements (Such as Jaccard, Jaro-Winkler, Sorensen) & study-learning methods (LRAR). In experimental results, Jaro-Winkler lends great efficiency compared to observations for LRAR, which introduces top another machine learning tests for the client. In together cases, the performance recommendations are so good & the developer framework helps the clients with altered regular exploration machine learning challenges [11].

Vaishali Bajpai et al. [2018] Introduce Dynamic Recommendation Systems to provide recommendations about the customer's interest. In this dynamic recommendation system, the use of web usage information is first used to detect user behavior and to measure similar user behavior score. In the second step, the product information will be collected based on the current search information about the customer who will assess the Sentiment Score and the Social Media Popular Score. The other uses qualitative matrix on the side of the customer's purchasing power. These three components- The Same Score behavior, Sentimental Score and Popular Score are used to calculate the combined weight for a particular product. Then a qualitative matrix and a computational weight are used to measure possible user recommendations. [12].

Dr. P.V.R.D. Prasada Rao et al. [2018] described different recommendation techniques. Many clustering, classification, association algorithms are used in building recommendation systems, rather than using individual algorithms if a combination of algorithms (hybrid approach) is used, recommendations can be predicted better [13].

Ouafae El Aissaoui [2018] Introduce an automated approach to finding web-based mining based on students learning style. Using a K-means algorithm, a definite learning style model (Silverman model, Felder) will classify students' log files. To check the efficacy of our activity, we utilize an actual world system collected after the E-Learning Structure. Our approach to providing promotional results shows experimental results [14].

Rahul  Kataria et al. [2017] Hybrid is applied to k-way and hydrogen for the Movielens database to get better film recommendations in this article. We restrained the presentation of our method to MAE, RMSE, SD, t-value. The experimental results of the Movielens Destination have specified that the discussion we discoursed has provided for scrutiny and high levels of confidence and a consistent & modified movie recommendation system using a certain number of clusters. Evaluation matrix (a fixed number for clusters) is lower than other methods. If the initial partition does not function properly, there are a few limitations suggested that you can reduce efficiency. Future searches can be used by other natural species alternatives instead. The best-promised methodology algos utilized for clustering & optimization deliver better performance of speed & accuracy [35].

Ahmed I. Saleh et al. [2015] This contributes to the contribution of classifying methods to improve the functioning of the RSS. An Intelligence Adaptive Vertical Recommendation (IAVR) structure will be presented in this paper. IAVR mentions text documents associated with a particular domain. Mainly, the paper focusses on the $1^{st}$ stage of IAVR consisting of 2 elements. The second is a distiller, the $2^{nd}$ is a multi-class division. The specified distiller is used like a binary wired to select the documents associated with the desired domain. A neuro-fuzzy system and a K Nearest Neighbour (KNN) are categorized. On the other hand, a multi-layer woven bees (NB) classifier merge a new system, depending on the specific learning technique called Association Study with Association rules. The efficiency of classifiers recommended by experimental results has proven, and everything promotes the system's accuracy [36].

M. Karthik [2013] The use of web-mining technology in e-commerce websites provides e-commerce sites with security. The client performance on the web is built on web mining, e-commerce & security. Consumer behavioral design analyzes e-commerce websites to improve. Various web mining algorithms & security algorithms deliver safety on e-commerce websites. Web mining techniques for example PageRank & Trade rank are utilized to improve web mining framework on an e-commerce website. Usually, web mining is built on web content or web usage mining. Web mining in this suggested system includes web structure mining, web content mining, decision analysis, & security analysis [37].

## VIII. CONCLUSION

This paper has a survey about the recommendation system which is the most significant field of web usage mining. There are some problems exist in the recommendation system. These can be solved using different classification techniques like K-nearest neighbor (KNN), Neural Networks (NN), Support Vector Machines (SVM), Bayesian classifier, Decision Tree, Rocchio's. But in all, the K-NN algorithm is the best classifier of classification to recommending online and real-time website to classify visitors/clients clickstream information. Actually Simple Syndication (RSS) reader website utilized to delivering applicable data towards the separate deprived of clearly examining aimed at it.

## REFERENCES

[1] Pradnya P. Soundwave, "Overview of Predictive and Descriptive Data Mining Techniques" IJARCSSE, Volume 5, Issue 4, April 2015.

[2] O. Etzioni. The World-Wide Web: Quagmire or Gold Mine? Communications of the ACM, 39(11):65–68, 1996.

[3] Adeniyi, D. A., Wei, Z., & Yongquan, Y. (2016). Automated web usage data mining and recommendation system using K-Nearest Neighbor (KNN) classification method. Applied Computing and Informatics, 12(1), 90–108.

[4] K.Amutha and Dr.M.Devapriya, "Web Mining: A Survey Paper", International Journal of Computer Trends and Technology (IJCTT) – Volume 4 Issue 9– Sep 2013, pp. 3038-3042.

[5] Prakash S Raghavendra et .al: Web Usage Mining using Statistical Classifiers and Fuzzy Artificial Neural Networks Infonomics Society 2011.

[6] Poonam Chavan, "Analytical study on Collaborative Filtering techniques for Location-based Recommendation", International Journal of Scientific & Engineering Research, Volume 6, Issue 8, August-2015, pp. 1750-1758.

[7] Teh, P. L., Ghani, A. A. A., & Chan Yu Huang. (2008). Survey on application tools of Really Simple Syndication (RSS): A case study at Klang Valley. 2008 International Symposium on Information Technology, pp. 1-8.

[8] Ronald J. Glotzbach, James L. Mohler, and Jaime E. Radwan, "Really Simple Syndication (RSS): An Educational Approach", JIME http://jime.open.ac.uk/2009/03.

[9] Upendra Singh and Saqib Hasan, "Survey Paper on Document Classification and Classifiers", International Journal of Computer Science Trends and Technology (IJCST) – Volume 3 Issue 2, Mar-Apr 2015, pp. 83-87.

[10] Alka Lamba and Dharmender Kumar, "Survey on KNN and Its Variants", International Journal of Advanced Research in Computer and Communication Engineering Vol. 5, Issue 5, May 2016, pp. 430-435.

[11] Werneck, R. de O., de Almeida, W. R., Stein, B. V., Pazinato, D. V., Júnior, P. R. M., Penatti, O. A. B., … Torres, R. da S. (2018). Kuaa: A unified framework for design, deployment, execution, and recommendation of machine learning experiments. Future Generation Computer Systems, 78, 59–76.

[12] Vaishali Bajpai, Yagyapal Yadav, "Survey Paper on Dynamic Recommendation System for E-commerce", International Journal of Advanced Research in Computer Science, Volume 9, No. 1, January-February 2018.

[13] Dr. P.V.R.D. Prasada Rao, Vineetha B, Ch. Priyanka, V. Satish, "A Comparative Study on Recommendation System using Hybrid Approach", International Journal of Mechanical Engineering and Technology (IJMET) Volume 9, Issue 1, January 2018.

[14] Ouafae EL AISSAOUI et al "Integrating web usage mining for an automatic learner profile detection: A learning styles-based approach", 978-1-5386-4396-9/18/$31.00 ©2018 IEEE.

[15] Rahul Kataria "An effective collaborative movie recommender system with cuckoo Search" Egyptian Informatics Journal 18 (2017) 105–112.

[16] Saleh, A. I., El Desouky, A. I., & Ali, S. H. (2015). Promoting the performance of vertical recommendation systems by applying new classification techniques. Knowledge-Based Systems, 75, 192–223.

[17] M.Karthik and S.Swathi, "Secure web mining framework for e-commerce websites", ISSN: 2231-2803.