# A WEBPAGE RECOMMENDER SYSTEM FOR ANONYMOUS USERS BASED ON WEBPAGES CLUSTERING

**Dr. R. B. Wagh**, Associate Professor, Dept. of CSE (Data Science), R.C.P.I.T., Shirpur
Email: rajnikantw@gmail.com

**Abstract:**
A Content Based web Personalization model is tailoring a website to the demands of certain individual or group of individuals. For accurate prediction of human subsequent movements, it employs a variety of data mining methods including mining association rules, sequential pattern finding, clustering, classification, and so on. Current online personalisation approaches lag in accurately predicting user interests. We have created a revolutionary online customization approach to deliver excellent suggestions. The suggested approach is centred on determining suitable weights as one of a blog's web sites. With this goal, we employed the distance measure of the visiting relation as well as the appearance frequency measurement of web sites. The method employs an improved graph-based partition technique to cluster webpages and much more correctly characterise current user actions. The Limit is utilised to make a recommendations judgement amongst webpages. Our test findings show approximately 61% accuracy, 34% coverage, and 45% F1 measure. It aids in improving the customer's user experience.

**Keywords:** Web Personalization, Classification, Recommender Systems, Clustering, Web Usage Mining

## 1. INTRODUCTION

Web personalizing is the practise of providing method involves to website visitors. Instead of delivering a single, extensive experiences, web customization enables businesses to provide users with subjective experiences that are suited to their requirements and aspirations.

Every day, a great deal of data as well as information is posted to the internet. Numerous websites, like news portals or other everyday life websites, are routinely updated on an hourly or every day basis. This necessitates the separation of useful and unnecessary data. Nevertheless, most web architectures are complex and huge in length. This could mislead consumers by providing them with superfluous and clear information. It therefore results in the loss of precious user time. It would be far preferable to eliminate the overload of data and save the user's time exploring. As a result, it is necessary to foresee the user's demands in order to optimise the user experience and provide them with whatever they desire. The best possible solution is web personalisation or a recommendation system.

To fulfil the goal of "Providing consumers with the data they need," a Based on web system may overtly or implicitly request limited inputs from users [1].

Website personalizing is also referred to in the literature also as practice of offering beneficial connections, things, and things to the user so as to save important time. Such availability is dependent on either the customer's stated mentions or what the system learns intuitively. Visitors are assessed based on their surfing history, geographic region, comparable users, things, links, products, and so on. During monitor customer and suggestion, many data mining approaches are utilised. It aids in the improvement of different E-Commerce platforms' business or consumer happiness [2].

Recommendation systems are generally concerned with overall score sites in which users score the items of the site. The item might be a film, song, literature, humour, or other such product. It's also an issue in research area with several implications. Amazon.com, for instance, suggests Discs, Publications, Audio, and Videos, as well as other things. Headlines at VERSIFI technology videos by Cinema lenses, as well as numerous others are instances [3]. Second section of this article offers a brief introduction of the components, processes, and kinds of online personalisation. Third Section

addresses web usage mining efforts for web personalisation. Section four is about methods, while Fifth section is about system assessment including experimental findings.

## 2. RELATED WORK

Other variations include content-based, data aggregation, rule-based, and browsing mining. This part will address some current online customization research, with a focus upon Web use mining. A recommendation system known as WebPUM was created by Jalali M. et al. This uses the longest common sequence algorithm (LCS) to identify user movement patterns and forecast users' future intents. They put out a method for categorising user browsing patterns in order to forecast users' future goals in order to successfully deliver online predictions [5].

Web customization is the subject of a review by Malik Z. and Fyfe. The core components of online personalization-learning, matching, and suggestions thoroughly addressed along with current developments. The merits and drawbacks of the corresponding phase and its variants-content-based, data aggregation, rule-based learners, and metaheuristic discussed. The difficulties associated with big data scalability, poor performance, black box filtering, accurate suggestion, and privacy concerns present fresh possibilities for study. The relevance of online personalisation for usage in e-commerce websites is also discussed by the researchers [4].

Liu and Keselj described a method for automatically classifying online user browsing patterns in order to forecast users' next actions. The technique relies on extracting web application records and the content of fetched web pages simultaneously. The information of a website page is represented using character N-grams. It and the web application log are merged to represent user profiles. The method was evaluated in an exercise to see how well it could classify and predict data. They achieved predictions and accuracy rate of about 65% & 70%, respectively [6].

In order to determine the subjects that consumers are interested about, Yang et al. created a graph-based iterations technique, that is employed to calculate user similarity. The researchers designed a particular matter Markov framework for understanding users' browsing habits, which captures both of the temporal as well as topical significance of webpages [7] in order to propose conceptually coherent sites.

Classifications as well as pre - processing of Internet information, the identification of relationships between this data, and determining the actions that the system should do are some of its essential components. Each of the below is a type of web data, Users' profiles: It offers details on website visitors. Each user's name, age, sex, nationality, state, relationship status, profession, and interests are included. Moreover, it includes details about the interests and tastes of the user. Such data is gathered by surveys or voter registrations, or it may be produced through Web application log data [8].

Simple text, pictures, or structured data, such as facts gleaned from databases, can all be considered content data. The final user is shown the content data. Structure Data: This phrase describes how stuff is arranged. They might be hyperlinks linking one website to another or HTML tags (data elements used inside a Web page). Links between different pages are present in order to limit user navigation to predetermined routes. Data used: It displays the traffic to a website. It contains information such as the user's Internet address, visit time and date, full route (folders or files) viewed, and referrer's URL [9, 15].

The components of Web personalisation are as follows. It consists of categorising and preparing Online data, extracting relationships between these data, and recognising the system's behaviours. • User Profile Information: Everything just provides information about an online majority of users. It comprises information such as a user's name, age, gender, nation, area, family status, schooling, and interests. Also it contains data on the user's tastes and hobbies. Such data is gathered by surveys or registration forms, or it might be created by analysing Web syslog messages [10].

User sessions that defined by Perkowitz and O. Etzioni et. al. created an Adapted Web site segmentation graph based technique. Such websites may automatically enhance the arrangement and display again for end customer. They mined knowledge from usage logs. With the Page Gather technique, a novel clustering strategy is given. Clusters are groups of cliques or components that are

coherently related. Larger clusters are easier and faster to compute on. The Page Gather method generates an index page containing hyperlinks to every one of the sites in a cluster [14, 17].

We used our findings to the DePaul University CTI event log collection (www.cs.depaul.edu). The data for the primary DePaul CTI Web server is contained in this data collection. The information is calculated using a random sampling of visitors that visited this site for two weeks in April 2002. Figure 1 visualize the raw (unfiltered) data set comprised 20950 visits from 5446 users. The filtering data files are created by screening low support page visits and removing session with a size of one. There are 13745 session and 683 page hits in the filtered data.
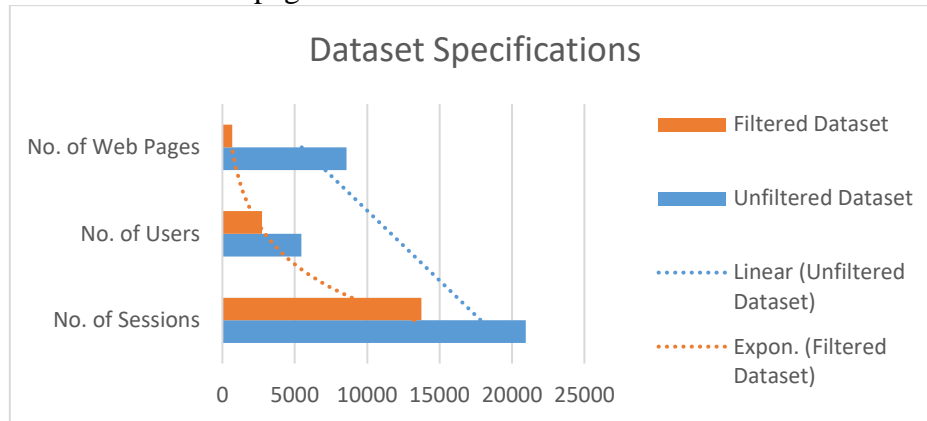


Figure 1: Original Dataset Specifications

Two experiments are carried out utilising the implementation specs described above. After modelling Web pages, the DFS method is used to cluster Web pages during the initial experiment. The second study involves a user established connection.

A recommendation systems is one that attempts to anticipate or filter priorities based on the user's selections. Recommendation systems were applied in a range of domains like films, audio, media, publications, scholarly articles, web searches, community tagging, and items in broad.

The Kernel generates choices based on

Suggestions Regarding the Content
Suggestions Regarding Similarity
Suggestions Depending on Collaborative Screening
Suggestions Depending on Surprising Library Factors
Factors affecting the development upon Embedding

## 3. CONTRIBUTIONS TO RESEARCH

The following seem to be the planned research work's unique contributions:

- Create new connection measures according to the distance between new website requests in session and the occurrence frequency to get an adequate relationship between web sites. These have not yet been built or suggested in either web page recommender systems.
- The recommended evaluation attempt to offer equal weight including all website pages. It takes into account every one of the session (13745 in total) when both website content were present. The measurements are constructed in such a manner that they provide normalised values ranging from zero to one.
- Creation of a virtualized graph that corresponds to the connection matrix. This graph's node are indeed the internet sites of such a website.
- Using the improved depth-first selection method, divide a cluster of networks. Such generated clusters are merely representations of previous users' browsing patterns. To identify more efficient navigation strategies, we employed edge thresholds and also cluster size thresholds during partitioning.
- Suggestions for the other web pages within this cluster. For this purpose, we firstly ordered the web sites by their final weighting factor, then imposed threshold levels and suggested just

those web content that met the threshold requirements. The suggested approach is evaluated using several performance metrics such as Visiting Coherent, Amount of Outliers, Cluster Number, Accuracy, Covering, and F1 upon that CTI database.

## 4. PERFORMANCE EVALUATION

A variety of strategies are utilized to assess online personalizing and recommendation systems. Our web customization platform is assessed using the different indicators, which are used both online as well as offline.

Visit Coherence: This measure is used to evaluate the level of quality of clusters created during the retrieval phase. It determines the proportion for Web pages within a user's account that correspond to a clusters representing the session under consideration [16].

We assessed visit coherence via dividing the information set into two halves: testing and training. The clustering procedure is utilised for the initial half of a data, while the remaining half serves to evaluate the level of quality of clustering created from the trained portion. This parameter is specified as the amount of internet pages in each session.

$$\beta_i = \frac{|\{p \in Si| \, p \in ci\}|}{N_i}$$

Here p represents a page, Si represents the ith sessions, Ci represents the cluster reflecting session I and Ni indicates the amount of page in the ith session. The average value for β for all N session in the dataset's assessment section is presented as:

$$\alpha = \frac{\sum_{i=1}^{n} \beta_i}{N_S}$$

Here is the proportion of visit-coherence that must be evaluated for the edge Threshold ranges [16].

Outliers: An outlier is indeed a percentage among Web sites that don't fit into another navigational patterns (cluster) so they do not contribute towards the online phase. Outliers are computed using edge Threshold settings. We attempted to reduce the number of outliers. So more exceptions there are, the less effective the cluster is, and inversely.

Accuracy: Its metric is used to assess the performance of suggestions. It denotes the amount of pertinent Web pages obtained by the user divided by the total amount of Web pages in the recommended set. We begin by dividing the collection into two parts: testing and training. Every navigational pattern npi (a dataset session) inside the test dataset is broken into two sections. The very first n npi page visits

$$\text{Accuracy}\{\text{Rec.Set}_{(AS,ET)})\} = \frac{|\text{Rec.Set}_{(AS,ET)} \cap \text{Test Set}_{(np-n)}|}{\text{Rec.Set}_{(AS,ET)}}$$

Here Rec. Sets (AS, ET) is the produced Recommended set in relation to a Active Sessions window (AS) as well as the Edge Threshold (ET). Tests Setnp-n is a test dataset session that discards its first n page visits. The similar Web pages that exist in both the recommended set as well as the test set.

F1: F1 is a Harmonic Mean of Accuracy and Coverage. It achieves maximum value when both Accuracy and Coverage achieve maximum values. It is given by

$$F1 = \frac{2 * \text{Accuracy}\{\text{Rec.Set}_{(AS,ET)}\} * \text{Coverage}\{\text{Rec.Set}_{(AS,ET)}\}}{\text{Accuracy}\{\text{Rec.Set}_{(AS,ET)}\} + \text{Coverage}\{\text{Rec.Set}_{(AS,ET)}\}}$$

Table 1: Performance evaluation (a) TC Matrix Produced by Previous Work, (b) Distance Matrix Produced by Proposed Work, (c) FC Matrix Produced by Previous Work, (d) Occurrence Frequency Matrix Produced by Proposed Work, (e) Weight Matrix Produced by Previous Work, (f) Relationship Matrix Produced by Proposed Work

**TC Matrix Produced by Previous Work**

|  | P1 | P2 | P3 | P4 | P5 | P6 | P7 | P8 | P9 |
|---|---|---|---|---|---|---|---|---|---|
| P1 | 0 | 0.25 | 0.55 | 0.38 | 0.54 | 0.35 | 0.18 | 0.4 | 0.63 |
| P2 | 0.25 | 0 | 0.37 | 0.18 | 0.27 | 0.18 | 0.38 | 0.34 | 0.35 |
| P3 | 0.55 | 0.37 | 0 | 0.18 | 0 | 0 | 0.74 | 0 | 0 |
| P4 | 0.38 | 0.18 | 0.18 | 0 | 0.15 | 0.34 | 0 | 0.56 | 0 |
| P5 | 0.54 | 0.27 | 0 | 0.15 | 0 | 0 | 0 | 0 | 0 |
| P6 | 0.35 | 0.18 | 0 | 0.34 | 0 | 0 | 0.17 | 0.28 | 0.29 |
| P7 | 0.18 | 0.38 | 0.74 | 0 | 0 | 0.17 | 0 | 0.31 | 0.56 |
| P8 | 0.4 | 0.34 | 0 | 0.56 | 0 | 0.28 | 0.31 | 0 | 0.18 |
| P9 | 0.63 | 0.35 | 0 | 0 | 0 | 0.29 | 0.56 | 0.18 | 0 |

(a)

**Distance Matrix Produced by Proposed Work**

|  | P1 | P2 | P3 | P4 | P5 | P6 | P7 | P8 | P9 |
|---|---|---|---|---|---|---|---|---|---|
| P1 | 0 | 0.3 | 0.6 | 0.4 | 0.6 | 0.45 | 0.2 | 0.42 | 0.7 |
| P2 | 0.3 | 0 | 0.4 | 0.2 | 0.3 | 0.2 | 0.4 | 0.4 | 0.4 |
| P3 | 0.6 | 0.4 | 0 | 0.2 | 0 | 0 | 0.8 | 0 | 0 |
| P4 | 0.4 | 0.2 | 0.2 | 0 | 0.2 | 0.4 | 0 | 0.6 | 0 |
| P5 | 0.6 | 0.3 | 0 | 0.2 | 0 | 0 | 0 | 0 | 0 |
| P6 | 0.45 | 0.2 | 0 | 0.4 | 0 | 0 | 0.2 | 0.3 | 0.35 |
| P7 | 0.2 | 0.4 | 0.8 | 0 | 0 | 0.2 | 0 | 0.4 | 0.6 |
| P8 | 0.42 | 0.4 | 0 | 0.6 | 0 | 0.3 | 0.4 | 0 | 0.22 |
| P9 | 0.7 | 0.4 | 0 | 0 | 0 | 0.35 | 0.6 | 0.22 | 0 |

(b)

**FC Matrix Produced by Previous Work**

|  | P1 | P2 | P3 | P4 | P5 | P6 | P7 | P8 | P9 |
|---|---|---|---|---|---|---|---|---|---|
| P1 | 0 | 0.51 | 0.34 | 0.32 | 0.29 | 0.75 | 0.63 | 0.51 | 0.78 |
| P2 | 0.51 | 0 | 0.44 | 0.73 | 0.73 | 0.52 | 0.37 | 0.3 | 0.3 |
| P3 | 0.34 | 0.44 | 0 | 0.56 | 0 | 0 | 0.63 | 0 | 0 |
| P4 | 0.32 | 0.73 | 0.56 | 0 | 0.45 | 0.3 | 0 | 0.35 | 0 |
| P5 | 0.29 | 0.73 | 0 | 0.45 | 0 | 0 | 0 | 0 | 0 |
| P6 | 0.75 | 0.52 | 0 | 0.3 | 0 | 0 | 0.29 | 0.82 | 0.75 |
| P7 | 0.63 | 0.37 | 0.63 | 0 | 0 | 0.29 | 0 | 0.34 | 0.32 |
| P8 | 0.51 | 0.3 | 0 | 0.35 | 0 | 0.82 | 0.34 | 0 | 0.61 |
| P9 | 0.78 | 0.3 | 0 | 0 | 0 | 0.75 | 0.32 | 0.61 | 0 |

(c)

**Occurrence Frequency Matrix Produced by Proposed Work**

|  | P1 | P2 | P3 | P4 | P5 | P6 | P7 | P8 | P9 |
|---|---|---|---|---|---|---|---|---|---|
| P1 | 0 | 0.57 | 0.4 | 0.33 | 0.33 | 0.75 | 0.66 | 0.57 | 0.85 |
| P2 | 0.57 | 0 | 0.5 | 0.8 | 0.8 | 0.57 | 0.4 | 0.33 | 0.33 |
| P3 | 0.4 | 0.5 | 0 | 0.66 | 0 | 0 | 0.66 | 0 | 0 |
| P4 | 0.33 | 0.8 | 0.66 | 0 | 0.5 | 0.33 | 0 | 0.4 | 0 |
| P5 | 0.33 | 0.8 | 0 | 0.5 | 0 | 0 | 0 | 0 | 0 |
| P6 | 0.75 | 0.57 | 0 | 0.33 | 0 | 0 | 0.33 | 0.85 | 0.85 |
| P7 | 0.66 | 0.4 | 0.66 | 0 | 0 | 0.33 | 0 | 0.4 | 0.4 |
| P8 | 0.57 | 0.33 | 0 | 0.4 | 0 | 0.85 | 0.4 | 0 | 0.66 |
| P9 | 0.85 | 0.33 | 0 | 0 | 0 | 0.85 | 0.4 | 0.66 | 0 |

(d)

**Weight Matrix Produced by Previous Work**

|  | P1 | P2 | P3 | P4 | P5 | P6 | P7 | P8 | P9 |
|---|---|---|---|---|---|---|---|---|---|
| P1 | 0 | 0.34 | 0.42 | 0.35 | 0.38 | 0.48 | 0.28 | 0.45 | 0.7 |
| P2 | 0.34 | 0 | 0.4 | 0.29 | 0.39 | 0.27 | 0.37 | 0.32 | 0.32 |
| P3 | 0.42 | 0.4 | 0 | 0.27 | 0 | 0 | 0.68 | 0 | 0 |
| P4 | 0.35 | 0.29 | 0.27 | 0 | 0.23 | 0.32 | 0 | 0.43 | 0 |
| P5 | 0.38 | 0.39 | 0 | 0.23 | 0 | 0 | 0 | 0 | 0 |
| P6 | 0.48 | 0.27 | 0 | 0.32 | 0 | 0 | 0.21 | 0.42 | 0.42 |
| P7 | 0.28 | 0.37 | 0.68 | 0 | 0 | 0.21 | 0 | 0.32 | 0.41 |
| P8 | 0.45 | 0.32 | 0 | 0.43 | 0 | 0.42 | 0.32 | 0 | 0.28 |
| P9 | 0.7 | 0.32 | 0 | 0 | 0 | 0.42 | 0.41 | 0.28 | 0 |

(e)

**Relationship Matrix Produced by Proposed Work**

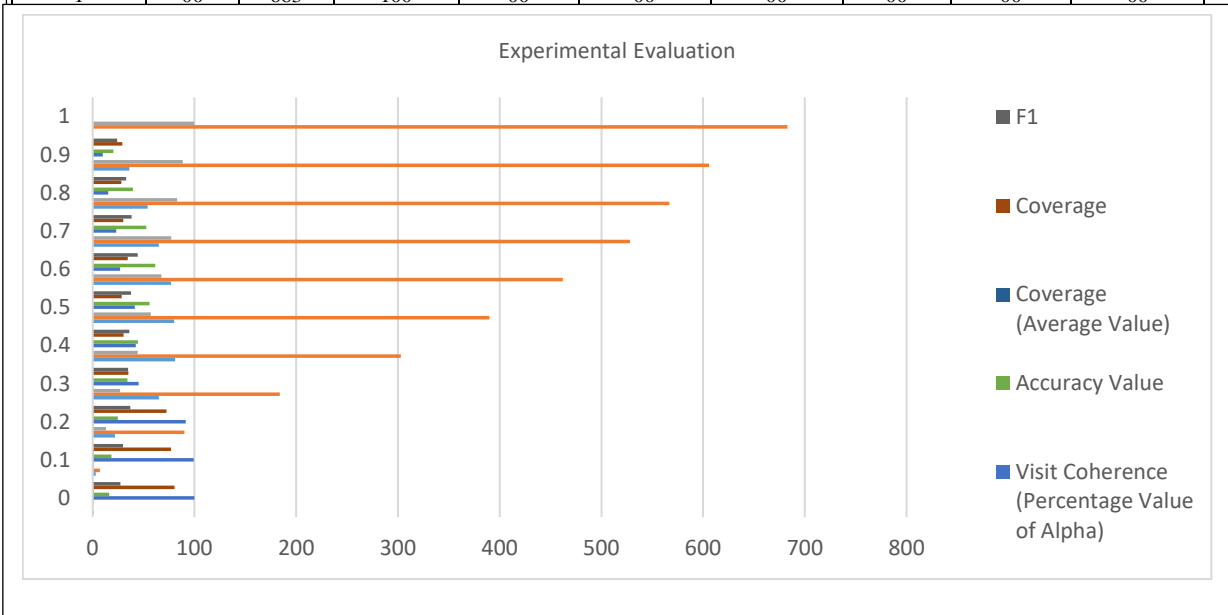|  | P1 | P2 | P3 | P4 | P5 | P6 | P7 | P8 | P9 |
|---|---|---|---|---|---|---|---|---|---|
| P1 | 0 | 0.39 | 0.48 | 0.36 | 0.42 | 0.56 | 0.3 | 0.48 | 0.76 |
| P2 | 0.39 | 0 | 0.44 | 0.32 | 0.43 | 0.29 | 0.4 | 0.36 | 0.36 |
| P3 | 0.48 | 0.44 | 0 | 0.3 | 0 | 0 | 0.73 | 0 | 0 |
| P4 | 0.36 | 0.32 | 0.3 | 0 | 0.28 | 0.36 | 0 | 0.48 | 0 |
| P5 | 0.42 | 0.43 | 0 | 0.28 | 0 | 0 | 0 | 0 | 0 |
| P6 | 0.56 | 0.29 | 0 | 0.36 | 0 | 0 | 0.24 | 0.44 | 0.49 |
| P7 | 0.3 | 0.4 | 0.73 | 0 | 0 | 0.24 | 0 | 0.4 | 0.48 |
| P8 | 0.48 | 0.36 | 0 | 0.48 | 0 | 0.44 | 0.4 | 0 | 0.33 |
| P9 | 0.76 | 0.36 | 0 | 0 | 0 | 0.49 | 0.48 | 0.33 | 0 |

(f)

Figure 2 : Visualization of Performance evaluation (a) TC Matrix Produced by Previous Work, (b) Distance Matrix Produced by Proposed Work, (c) FC Matrix Produced by Previous Work, (d) Occurrence Frequency Matrix Produced by Proposed Work, (e) Weight Matrix Produced by Previous Work, (f) Relationship Matrix Produced by Proposed Work

The raw data set comprised 20950 visits from 5446 individuals. The filter information files are generated by screening weak support page visits and removing visits with a size of one. Table 3 specifies there are 13745 session & 683 page hits inside the input data. Two tests are carried out utilising the implementation specs described previously. After modelling Web pages, the DFS method is used to group Web pages during the initial experiment. In the second study, the LCS method is used to classify a user's current session into one of the clusters.

Table 1: Experimental Results

| Edge Threshold Value | No. of Clusters | No. of Outlier Web Pages | Percentage of Outliers | Visit Coherence (Average Value of Alpha) | Visit Coherence (Percentage Value of Alpha) | Accuracy Value | Coverage (Average Value) | Coverage (Percentage Value) | F1 |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 01 | 00 | 00 | 01 | 100 | 16.33 | 0.80 | 80.40 | 27.14 |
| 0.1 | 03 | 07 | 1.02 | 0.99 | 99.28 | 18.42 | 0.77 | 77.01 | 29.72 |
| 0.2 | 22 | 90 | 13.17 | 0.91 | 91.49 | 24.78 | 0.72 | 72.48 | 36.93 |
| 0.3 | 65 | 184 | 26.93 | 0.45 | 45.27 | 34.12 | 0.35 | 35.19 | 34.64 |
| 0.4 | 81 | 303 | 44.36 | 0.42 | 42.32 | 44.61 | 0.30 | 30.37 | 36.13 |
| 0.5 | 80 | 390 | 57.10 | 0.41 | 41.27 | 55.70 | 0.28 | 28.47 | 37.68 |
| 0.6 | 77 | 462 | 67.64 | 0.27 | 26.99 | 61.42 | 0.34 | 34.56 | 44.23 |
| 0.7 | 65 | 528 | 77.30 | 0.23 | 23.16 | 52.86 | 0.30 | 30.10 | 38.35 |
| 0.8 | 54 | 567 | 83.01 | 0.15 | 15.36 | 39.47 | 0.28 | 28.22 | 32.91 |
| 0.9 | 36 | 606 | 88.72 | 0.10 | 9.90 | 20.32 | 0.29 | 29.16 | 23.95 |
| 1 | 00 | 683 | 100 | 00 | 00 | 00 | 00 | 00 | 00 |



Experimental Evaluation

## 5. CONCLUSIONS AND FUTURE DIRECTIONS

We invented a new method of web customisation. The suggested approach is centred on determining more suitable weights as one of a website's web pages. A unique formula calculating the distance connection as well as the meet common is used to model the Web content. This association matrix's improved clustering assisted us in forming more appropriate groupings. We used LCS to categorise active users. The Threshold we employed in the last round of suggestions enhances our system's accuracy. Yet, it has an impact on coverage. Semantic comprehension of the fundamental domain may boost suggestion quality even further.

REFERENCES

[1] M. Eirinaki and M. Vazirgiannis, (2003), "Web Mining for Web Personalization", In Proceedings of ACM Transactions on Internet Technology (TOIT).ACM, Athens, Greece, 3(1), pp.1-38, http://doi.acm.org/10.1145/643477.643478

[2] G. Adomavicius and A. Tuzhilin, (2005), "Toward the Next Generation of Recommender Systems: A Survey of the State-of-the-Art and Possible Extensions", IEEE Transactions on Knowledge and Data Engineering, 17(6), pp.734–749.

[3] U. Gulden and M. Matthew, (2008), "Personalization Techniques and Recommender Systems: Series in Machine Perception and Artificial Intelligence", World Scientific Press, Vol. 70, Singapore,

[4] M. Bamshed and A. Sarabjot Singh, (2005), "Intelligent Techniques for Web Personalization", Springer, New York.

[5] R. Francesco, R. Lior, and S. Bracha, (2015), "Recommender Systems Handbook", Springer, New York.

[6] Z. Malik and C. Fyfe, (2012), "Review of Web Personalization", Journal of Engineering Technologies in Web Intelligence, 4(3).

[7] Q. Yang, J. Fan, J. Wang, and L. Zhou, (2010), "Personalizing Web Page Recommendation via Collaborative Filtering and Topic-Aware Markov Model", IEEE International Conference on Data Mining, 1(1), pp. 1145-1150.

[8] Y. AlMurtadha, N. Sulaiman, N. Mustapha and N. Udzir, (2011), "IPACT: Improved Web Page Recommendation System Using Profile Aggregation Based On Clustering of Transactions", American Journal of Applied Sciences, 8 (3), pp. 277-283.

[9] M. Jalali, N. Mustapha, N. Sulaiman and A. Mamat, (2010), "WebPUM: A Web-based Recommendation System to Predict User Future Movements", Expert Systems Applications, 37, pp. 6201-6212.

[10] H. Liu and V. Keselj,(2007), "Combined Mining of Web Server Logs and Web

[11] Contents for Classifying User Navigation Patterns and Predicting Users' Future Request", Data and Knowledge Engineering, 61(2), pp.304-330.Conference Short Name:WOODSTOCK'18

[12] B. Mobasher, R. Colley, and J. Shrivastav, (2000), "Automatic Personalization based on Web Usage Mining", Communications of the ACM, 43 (8), pp. 142-151.

[13] C. Sumathi, R. Valli and T. Santahnam,(2010),  "An Application of Session Based Clustering to Analyze Web Pages of User Interest from Web Log Files", Journal of Computer Science, 6(1), pp.785-793.

[14] Muhan Zhang , Yixin Chen,(2020), "Inductive Matrix Completion based on graph neural network", International Conference on ICLR,pp. 1-14.

[15] Stefen Rendle, Li Zhang, (2019), "On the difficulty of evaluating baselines: A Study on recommender systems", Google Research, Mountain View, pp. 1-19.

[16] Daeryong Kim, Bongwon Suh , (2019), "Enhancing VAEs for Collaborative Filtering: Flexible Priors & Gating Mechanisms", RecSys '19, September 16-20, 2019, Copenhagen, Denmark.pp.1-5.

[17] Vojtech Vancura, Pavel Kordık, "Deep Variational Autoencoder with Shallow Parallel Path for Top-N Recommendation (VASP)", (2021),pp.1-6.