# Predicting implicit search behaviors using log analysis

**L. Leema Priyadharshini [1] *, S. Florence [1], K. Prema [1], C. Shyamala Kumari [1]**

[1] *Assistant Professor, Department of Computer Science and Engineering, School of Computing, Vel Tech Rangarajan Dr. Sagunthala R&D Institute of Science and Technology, Avadi, Chennai-62, TamilNadu, India*
*Corresponding author E-mail: leemapriyadharshinil@veltechuniv.edu.in*

## Abstract

Search engines provide ranked information based on the query given by the user. Understanding user search behavior is an important task for satisfaction of the users with the needed information. Understanding user search behaviors and recommending more information or more sites to the user is an emerging task. The work is based on the queries given by the user, the amount of time the user spending on the particular page, the number of clicks done by the user particular URL. These details will be available in the dataset of web search log. The web search log is nothing but the log which contains the user searching activities and other details like machine ID, browser ID, timestamp, query given by the user, URL accessed etc., four things considered as the important: 1) Extraction of tasks from the sequence of queries given by the user 2) suggesting some similar query to the user 3) ranking URLs based on the implicit user behaviors 4) increasing web page utilities based on the implicit behaviors. For increasing the web page utility and ranking the URLs predicting implicit user behavior is a needed task. For each of these four things designing and implementation of some algorithms and techniques are needed to increase the efficiency and effectiveness.

*Keywords*: *Search Engine; Search Behavior; Tasks; Ranking; Implicit Behavior.*

## 1. Introduction

The user searching activities towards search engines records those information using web search logs. Previous papers shown that search logs used in various applications such as finding the satisfaction of the user, web page utilities of user, user interest towards the particular task, suggesting some useful queries, ranking the needed web pages while retrieving the information from the search engine, recommending some websites to the user etc., But almost all the previous papers work at the level of query (query trail) or at the level of session (session trail).

The web search log contains the following columns: 1. time 2. Event 3. Value 4. Task. The user can open multiple tabs at the same time for giving multiple tasks. Each tab has the separate unique identification number for the perfect retrieval of information. The previous papers analysis the user search behavior only based on the session level trails but this work focuses the task level analysis to predict the user search behaviors. Applying the data mining techniques and algorithms is very much necessary for the prediction of user behaviors.

Web search logs contain user activities towards search engines, such as queries, reading time and clicks. Search trails store the prints left by users of search engines in their search processes. In the previous literature of understanding search trails, most of the previous projects have applied in the work of user satisfaction prediction, ranking function identification, query suggestion, etc.

## 2. Related work

In 2014, Zhen Liao, Yang Song proposed the concept of grouping the user queries into relevant tasks for predicting the user behaviors like amount of clicks towards the URL and reading time of the particular webpage etc.,[31] They have implemented this based on the implicit behaviors of the user. They have implemented query suggestion, and URL recommendation for predicting user behaviors.

Heasoo Hwang and Hady W. Lauw, proposed the concept organizing the user search histories from the web search logs based on the queries given by the user. The queries will vary depending upon the information need of the user. They have grouped the users historical queries in a dynamic and automated fashion. They have shown that this approach is very useful in query suggestion, result ranking based on the clicks. They have identified the different approaches for these things [1].

### 2.1. User behavior prediction

The proposed technique for clustering queries into task shows the less efficiency. There are lot of ranking methods have been identified in the literatures. Zhicheng Dou and Ruihua Song proposed the method of evaluating the effectiveness personalized web search. They have implemented the techniques for incorporating the personalization by using the implicit behaviors of the users such as clicking rate and the dwell time of the user. They have taken the query logs of the windows live search for doing their experiments. They have predicted the performance of the personalization using some evaluation techniques [2].

Yufei Tao and Cheng Sheng proposed that finding the nearest neighbors using the keyword with the spatial data mining algorithms which will be useful for the prediction of user behaviors. Range search and nearest neighbor retrieval techniques are used to predict these things [22].

## 2.2. Ranking

Yongdong Zhang and Xiaopeng Yang proposed the techniques to boost up the text based retrieval of images through image re-ranking methods. They have implemented this based on the click based query dependent relevance feedback mechanisms. They have proposed the novel re-ranking algorithm for ranking images. They have considered the relevance feedback as an implicit thing [4]. John b. Killoran proposed the techniques for increasing the website visibility using the search engine optimization techniques. They have answered the different questions related to search engine ranking mechanisms and mechanisms to increase the visibility of web pages. This paper proposed the concept of increasing the web page visibility to the user using some optimization technique. They have proposed the click through rate and the bouncing rate.

GU Hong, ZHAO Guangzhou and QIU Jun proposed the online learning approach for ranking mechanisms. This paper deals with relevance feedback mechanism which sends the user search behavior automatically [24].

## 2.3. Utility identification

A.K. Shanna and Neha Aggarwal proposed the web based search result optimization technique for mining of web search logs contains the queries given by the user. The optimization technique is based on the historical query logs which predict the information need of the user. It will reduce the navigation time of the user result set. They have applied the clustering technique in their query logs and capture the particular pattern of clicking of web pages while searching the information [8]. Kinam Park and Sangyep Nam proposed the techniques to extract the search intentions from the web search log of the particular search engine. They have built the intention graph to understand their information needs. It extracts the user search intentions using the clustering algorithms and also the labeling algorithms. The extracted user intentions are represented in the particular intention graph to identify the user satisfaction scores. They have conducted the evaluation measuring to measure the effectiveness of the algorithm [9]. Christos Zaroli-agis and Athanasios Papagelis proposed the approach for increasing the web search done by the user. They have analyzed the web page usages of the user like number of documents the user shared to others, number of pages the user downloaded etc., Finally they have done the page rank mechanism to rank the URLs based on the in-links and the out-links[13].

## 2.4. Task clustering

A task of the user can be seen as a set of meaningfully relevant search query trails within single session. Since it use the beginning search query of each particular query trail to be represent the whole query trail, a task may be simply represented by all such beginning search queries from the particular user query trails. Dilek Hakkani-Tur and Gokhan Tur proposed the techniques to understand the task intention of the user or domain identification by using user search query logs and the click logs. Initially the natural language query is given to the syntax based transformation to convert the natural languages to the normal query terms and then the efficiency was calculated. The calculation shows that it improves significantly for efficient domain detection preferable in web based utterances [14]. Ryen W. White, Mikhail Bilenko and Silviu Cucerzan proposed the web search interaction mechanism by conducting the study. They have done the TRAIL extraction by grouping the interaction logs based on browser Id. They have done the work for predicting the destination information. They have conducted the study for baseline, query suggestion and query destination. Based on the previous searches done by the users they suggested the websites to the users who requesting the similar queries. Based on the study and the experimental results they improved the interaction between the user and the web [18].

## 2.5. Query suggestion

As a particular user enters a particular query in the searching box of particular search engine, a drop-down menu will appears with the query suggestions to complete the user query. Yang Cao and Ju Fan proposed the three effective techniques for error correction, query suggestion and query expansion for increasing the usability of the patent searching facility. These things help the user to increase the usability over patents towards the users who are searching for the effective patents.

## 3. Goal

Extract the tasks from the web search logs contains queries given by the user. Suggest the relevant queries based on the task. Predict the user behavior and increase the web search interaction.

### 3.1. Query clustering

A task of the user can be seen as a set of meaningfully relevant search query trails within single session. Since it use the beginning search query of each particular query trail to be represent the whole query trail, a task may be simply represented by all such beginning search queries from the particular user query trails. The problem remaining is to combine the similar user queries in together. However, how unambiguously (meaningfully) define the semantic relevance between the user queries remains very much challenging.

### 3.2. Query suggestion

The query suggestion is a service which provides query suggestions that can able to complete a user's searching query. As a particular user enters a particular query in the searching box of particular search engine, a drop-down menu will appears with the query suggestions to complete the user query.

### 3.3. User behavior prediction

This is to understand that whether a particular user was satisfied by the information given or not. After the search conducted by the user process, several implicit feedback techniques should be identified as measures for analyzing user satisfaction.

## 4. System architecture

This is to understand that whether a particular user was satisfied by the information given or not. After the search conducted by the user process, several implicit feedback techniques should be identified as measures for analyzing user satisfaction. The Figure 1 rep [resents the architectural diagram of the proposed work.
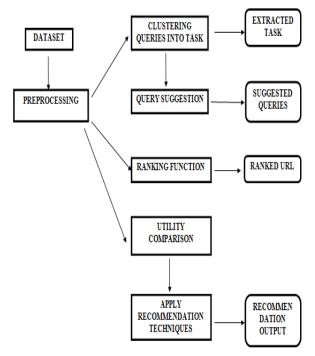
**Fig. 1:** Architectural Diagram.

## 4.1. Preprocessing

Preprocessing is the process of removing the stop word from the query given by particular user and makes the stemming in those words for obtaining the root words present in the particular word. A web search log contains a set of users, and each user has the sequence of consecutive activities.

This is the process of removing the stop word from the query given by particular user and makes the stemming in those words for obtaining the root words present in the particular word. For removing the stop word the stop word removal algorithm has been used.

## 4.2. Clustering queries into task

A web search log contains a set of users, and each user has the sequence of consecutive activities. A search activity is a single user query submitted to a particular search engine. In web search logs, a single query is followed by the sequence of browsing behaviors before the next user query is submitted by the particular same user. Thus, the simplest searching logs extraction is to treat one query and its followers as an free query trail.

## 4.3. Query suggestion

The query suggestion is a service which provides query suggestions that can able to complete a user's searching query. As a particular user enters a particular query in the searching box of particular search engine, a drop-down menu will appears with the query suggestions to complete the user query.

## 4.4. Ranking function

This is the process which takes the preprocessed dataset as input and based on the implicit behaviors such as click rate, it will rank the URLs for the particular task. This process fully based on the recently clicked URLs. Clicking rate is defined as the ratio between the number of sessions with clicks to the displayed URL for the given query and the total number of tasks in which the particular URL is viewed for the particular query.

## 4.5. Utility comparison

This is to understand that whether a particular user was satisfied by the information given or not. After the search conducted by the user process, several implicit feedback techniques should be identified as measures for analyzing user satisfaction.

## 5. Algorithms

Algorithm deals with clustering the queries into tasks. Initially queries are clustered based on the id. If multiple rows contain the same id then it will be grouped into the single id and queries are listed under that id. If the particular word is similar to the word of the query then the entire query will be given as suggestion for the particular word. Then print the words parallel to the word which needs suggestion.

Input: Preprocessed dataset P
Output: Clustered set of tasks $T_i$
Initialization: row=0
For len=1: Q-1 do
For j=0: len [Q] do
If $Q_{row}$ [len]$_j$= $Q_{row}$[len+1]$_j$ then //compare with all the words under particular id
$T_i$= $Q_{row}$ [len]$_j$ //Display words if similar words existing
Else
$T_i$= $Q_{row}$ [len] //Display query as task if not
For i=0: len[D] do
Words[i] =$Q_i$ // for extracting the words
If words[i] == $Q_i$[word] then // for checking sim of word and the word[query]
$Sug_i$=Display [$Q_i$] //copying query into the suggestion
Print words[i], $Sug_i$ //printing the suggestion parallel to the query
Return Ti

## 6. Experimental results

The implementation of the user behaviour prediction system builds upon the platform of java uses a integrated platform called netbeans. Java was chosen as the basis for this implementation due to its interoperable characteristic, which can call Google API and use the data present in the excel files. The extensible implementation includes ranking the URLs and the websites suggestion is also based on the language of java and the dataset of sogou search engine. The implementation is based on the dataset of sogou search engine and the algorithms of user behavior prediction.

### 6.1. Query suggestion

The implementation of query suggestion considered the dataset of search engine as input, which has above 1000 records containing the attributes like user id, query, click rate, URL clicked and timestamp. The implementation has the suggestion algorithm to implement the query suggestion process which takes the preprocessed dataset as input. The implementation takes each word of the query into an independent array and each array of words is considered as the input for the query suggestion, which then scans all the words present in the array of query. The suggested words will be displayed after the scanning and comparing process.

The figure 2 represents the implementation of query suggestion for some of the words of the user given queries.

**Fig. 2:** Query Suggestion.

## 6.2. Clustering queries into task

The implementation of task clustering takes the preprocessed dataset as input. It compares the words of the query given by user for clustering the queries into task. Based on the similar words and the different words it will be vary. It clusters the queries using the comparison process done in array of words present in the different queries issued by the users.

The figure 3 represents the implementation of task clustering.


**Fig. 3:** Clustered Task.

## 6.3. Ranking

The implementation of ranking URL is based on the number of clicks done by the user towards particular URL. Maximum number of clicks makes the URL to be ranked first. The ranking of the URL considered the task of the user to be an important one. The figure 4 represents the implementation of ranking in an offline mode.


**Fig. 4:** Ranked URL.

## 6.4. URL recommendation

The URL recommendation done by comparing the task of the user and the URLs retrieved for the particular query. The figure 5 represents the URL recommendation.


**Fig. 5:** Recommended URL.

The following figure 6 shows the performance of the proposed system. The performance of the proposed system is good in terms of cpu utilization and memory.


**Fig. 6:** Performance Evaluation.

The figure 7 represents the comparison between the attributes present in the dataset.


**Fig. 7:** Attribute Comparison.

The figure 8 represents the ROC curve of the proposed work.

**Fig. 8:** ROC Curve.

## 7. Conclusion

In this work the automated clustering and the query suggestion system has been implemented which clusters the given user queries into tasks by using the query comparison algorithm and then word suggestion is retrieved by using the query suggestion algorithm. The preprocessing algorithm for preprocessing the query is also implemented to get the output with high performance. The word which needs suggestion is compared with the queries of the query dataset and suggested queries are displayed when it is similar to the word which needs suggestion.

The proposed work given the expected output of task clustering and the query suggestion. The task clustering efficiently produces the task for the given set of queries under the particular user identification. The query suggestion also gives the efficient output for the given list of words in an efficient manner. Rank the user visited websites and also the recommendation of URL for the user under particular id have been done in an effective manner.

## References

[1] Heasoo Hwang, Hady W. Lauw, and Alexandros Ntoulas," Organizing User Search Histories", in *IEEE Trans. Knowl. and Data Engg., vol. 24, no. 5*, pp. 912-925, 2012. https://doi.org/10.1109/TKDE.2010.251.

[2] Zhicheng Dou, Ruihua Song, Ji-Rong Wen, and Xiaojie Yuan, "Evaluating the Effectiveness of Personalized Web Search", in *IEEE Trans. Knowl. and Data Engg., Vol. 21, No. 8*, pp. 1178-1190, 2009. https://doi.org/10.1109/TKDE.2008.172.

[3] Yang Cao and Guoliang Li, "A User-Friendly Patent Search Paradigm", in *IEEE Trans. Knowl. and Data Engg., Vol. 25, No. 6*, pp. 1439-1443, 2013. https://doi.org/10.1109/TKDE.2012.63.

[4] Yongdong Zhang, Senior Member, IEEE, Xiaopeng Yang, and Tao Mei, Senior Member, IEEE," Image Search Reranking With Query Dependent Click-Based Relevance Feedback", in *IEEE Trans. Image Processing, Vol. 23, No. 10*, pp. 4448-4459, 2014. https://doi.org/10.1109/TIP.2014.2346991.

[5] Christoph Kofler, Linjun Yang, Member, IEEE, Martha Larson, Member, IEEE, TaoMei, Senior Member, IEEE, Alan Hanjalic, Senior Member, IEEE, and Shipeng Li, Fellow, IEEE," Predicting Failing Queries in Video Search", in *IEEE Trans. Multimedia, Vol. 16, No. 7*, pp. 1973-1985, 2014. https://doi.org/10.1109/TMM.2014.2347937.

[6] John B. Killoran," How to Use Search Engine Optimization Techniques to Increase Website Visibility", in IEEE *transactions on professional communication, vol. 56, no. 1*, pp. 50-66, 2013. https://doi.org/10.1109/TPC.2012.2237255.

[7] Chen Gong, Keren Fu, Artur Loza, Qiang Wu, Member, *IEEE,* Jia Liu, and Jie Yang," PageRank Tracker: From Ranking to Tracking", in IEEE Trans. *Cybernetics, Vol. 44, No. 6*, pp. 882-893, 2014. https://doi.org/10.1109/TCYB.2013.2274516.

[8] A.K. Shanna, Neha Aggarwal, Neelam Duhan and Ranjna Gupta, " Web Search Result Optimization by Mining the Search Engine Query Logs", in IEEE *2010 Int. Conf. Meth. Models in Computer Science*, pp. 39-45, 2010.

[9] Kinam Park, Taemin Lee, Soonyoung Jung, Sangyep Nam, "Extracting Search Intentions from Web Search Logs", in *IEEE,* 2010.

[10] Jinjia Cheng, Chuanchang Liu, Yong Peng, "Expression Of User Personalized Search Behavior Based On Keyword Query Series And Bayesian Network", in *IEEE,* 2009.

[11] Omair Shafiq, Reda Alhajj, Jon G. Rokne, "Reducing Search Space for Web Service Ranking using Semantic Logs and Semantic FP-Tree based Association Rule Mining", in *Proc. IEEE 9th Int. Conf. Semantic Computing,* 2015. https://doi.org/10.1109/ICOSC.2015.7050771.

[12] Xueqing Gong, Xinyu Guo, Rong Zhang, Xiaofeng He and Aoying Zhou, " Search Behavior Based Latent Semantic User Segmentation for Advertising Targeting", in *IEEE 13th Int. Conf. Data Mining,* 2013. https://doi.org/10.1109/ICDM.2013.62.

[13] Athanasios Papagelis and Christos Zaroliagis, Member, IEEE, "A Collaborative Decentralized Approach to Web Search", in *IEEE Trans. Sys., Man, and Cybernetics, Vol. 42, No. 5*, pp. 1271-1290, 2012. https://doi.org/10.1109/TSMCA.2012.2187887.

[14] Dilek Hakkani-T`ur ,Gokhan Tur ,Larry Heck Asli, Celikyilmaz, Ashley Fidler Dustin Hillard ,Rukmini Iyer, Sarangarajan Parthasarathy, " Employing Web Search Query Click Logs for Multi-Domain Spoken Language Understanding", in *IEEE*, 2011.

[15] Shaoming CHEN, Yajun DU, Qiangqiang PENG, "Extracting query expansion terms based on user's search behavior", in *Sec. Int. Symp. Computational Intel. and Design, IEEE computer society*, 2009.

[16] Isak Taksa, Sarah Zelikovitz, Amanda Spink, "Using Web Search Logs to Identify Query Classification Terms", in *Int. Conf. Infor. Tech. IEEE*, 2007. https://doi.org/10.1109/ITNG.2007.202.

[17] Ricardo Baeza-Yates, Carlos Hurtado, Marcelo Mendoza and Georges Dupret, " Modeling User Search Behavior", in *Proc. of the Third Latin American Web Congress IEEE*, 2011.

[18] Ryen W. White, Mikhail Bilenko, Silviu Cucerzan,"Studying the Use of Popular Destinations to Enhance Web Search Interaction", in *Proc. ACM SIGIR*, 2007. https://doi.org/10.1145/1277741.1277771.

[19] Chien-Kang Huang, Lee-Feng Chien, Yen-Jen Oyang, " Relevant Term Suggestion in Interactive Web Search Based on Contextual Information in Query Session Logs", in *Journal American Soc. For Infor. Sci. And Tech., 54(7):638–649*, 2003.

[20] Rosie Jones, Benjamin Rey and Omid Madani, Wiley Greiner, "Generating Query Substitutions", in *ACM*, 2005.

[21] Craig Siverstein, Hannes marais, Monika Henzinger and Michael moraiz, "Analysis of very large web search engine query logs" , in *Compaq System Research*, 2000.

[22] Yufei Tao, "Fast nearest Neighbor Search with Keywords", in *IEEE Trans. Knowl. And Data Engg. Vol 25*, 2014.

[23] Chao Li, Bin Wu, "The Analysis of Youths' Searching Behavior", in *Int. Journal of Comp. And Comm. Engg,* 2011.

[24] GU Hong, ZHAO Guangzhou, "Online Metric Learning for Relevance Feedback in E-Commerce Image Retrieval Applications", in *Tsinghua Science And Tech.,Vol 16, No 4*, 2011.

[25] Mark Sanderson and W. Bruce Croft," The History of Information Retrieval Research", in *Proc. IEEE, Vol. 100*, 2012.

[26] Kai Li, "Running and Chasing - The Competition between Paid Search Marketing and Search Engine Optimization", in *47th Hawaii Int. Conf. System Science*, 2013.

[27] Fabrício Benevenuto, "Characterizing User Behavior in Online Social Networks", in *IMC'09*, 2009.

[28] Adish Singla," Studying Trailfinding Algorithms for Enhanced Web Search", in *SIGIR '10*, 2010.

[29] Yang Song, Hao Ma, Hongning Wang, Kuansan Wang," Exploring and Exploiting User Search Behavior on Mobile and Tablet Devices to Improve Search Relevance", in *ACM*, 2013.

[30] Laura A. Granka," Eye-Tracking Analysis of User Behavior in WWW Search", in *SIGIR '04*, 2004.

[31] Zhen Liao, Yang Song, Yalou Huang, Li-wei He, and Qi He, "Task Trail: An Effective Segmentation of User Search Behavior", in *IEEE Trans. Knowl. And Data Engg, Vol. 26, No. 12,* pp. 3090-3102, 2014. https://doi.org/10.1109/TKDE.2014.2316794.