



ENHANCED SEARCH RESULTS IN WEB AND IMAGE SEARCH ENGINE USING TASKTRAIL APPROACH

Ms.R.Nandini

PG Scholar, Dept of Computer Science
M.Kumarasamy College of Engineering, karur
nandini6.be@gmail.com

Mrs.V.Kavitha

Assistant Professor, Dept of Computer Science
M.Kumarasamy College of Engineering, karur
Kavithav.cse@mkce.ac.in

Abstract— In web search engine there are many techniques used for effective retrieval of information. Personalized web search is used to improve the quality of search services. When the user enters the query in search engine it provides large no of search results and it becomes important to rank these search results that it provides accurate and more relevant results. These are performed using ranking algorithm. Every user's search interest are different from other user. In existing system, the search engine that provides results for a given user query is not depend on their previous searching behavior. Here "task trail" is used to understand the user search behavior. Task is defined to an atomic user information need. Web search logs are mainly at session or query level where users can submit several queries within one task and handle several tasks within one session. The problem of task identification have been addressed by previous studies but a little advantage of using task over session and query for search applications. And implement personalization concept in web pages and web images in real time search engine environments

Keywords: Personalized web Search, Search engine, tasktrail, User Search Behavior

I. INTRODUCTION

Web search engines is an essential tool to address their duties. Suppose search engines return irrelevant results to the user means that does not satisfy their real intentions. For example, a real estate agent can issue the query "office" finding the vacant office space, but an IT specialist may issue the same query finding for popular Microsoft productivity software. Personalized web search is a search techniques to provide better search results. Many recent research efforts focused on this area. User information will be collected and analyzed the user intention behind the issued query.

In this module user enter the search query to retrieve the search results from the server. Then it analyze whether the query is personalized or general. Here online profiler is used to collect information from the user and then convert it into generalized profile. It includes search history and search experience. Greedy DP algorithm is used to improve the search results. The taxonomy repository is created for quick access. Prune leaf function is implemented to remove

unwanted search results. Greedy IL algorithm improves the efficiency using heuristics based method. Prune leaf operation is used to reduce the discriminating power.

There are two types in personalized web search click log based methods and profile based ones. The first method click log based methods to collect the user information via history. But it can only work on repeated queries [1]. The second method profile based method improve the complicated model. It is effective for almost all sorts of queries [1]. But this method collects the personal information implicitly from query history [2], [3], bookmarks pages [9], browsing history [5] [6], click-through data [7], [8], documents etc. Unfortunately, such implicitly collected personal data can easily affect the user's private life. It leads to a privacy problem [11]. To protect user privacy in profile based personalized web search have to consider two effects during the search process. First they try to improve the search quality with the utility personalization of the user profile. Second they try to protect the privacy contents existing in the user profile.

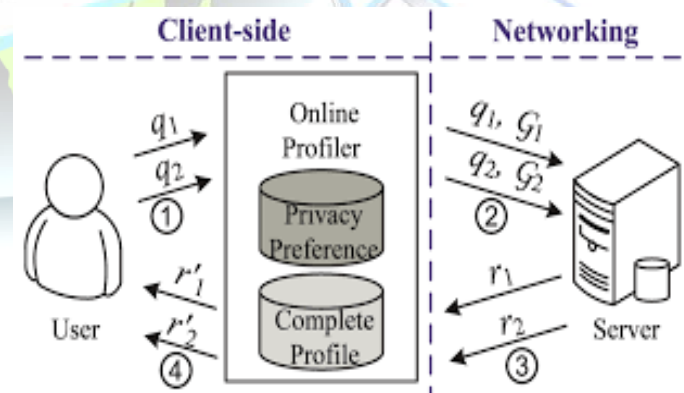


Fig 1. System architecture of UPS

UPS consist of search engine server and a number of clients. In personalization only a small portion of the user profile that is called generalized profile. The aim of personalized web search is user privacy can be protected without compromising the personalized search quality. The



privacy preserving personalized web search are far from optimal. A user profile is generalized only for offline and used to personalize all queries from the same user. The existing methods do not take the customization of privacy requirements. When creating personalized search results many personalization techniques require iterative user interactions.

II. RELATED WORK

A. Statistical language modeling analysis:

Long term search history includes rich information about a user's search preferences. Here statistical language modeling based methods to mine information from long term search history and it is used to estimate of the query language model. Using click through data of past searches that are related to the current query so it gives best performance. The decision of documents to return is based only on the query without consideration of a particular user's preferences. Suppose when a query is ambiguous the search results are mixed in content which is certainly non-optimal for the user.

There is more search contexts from the user's interests and needs, personal document collection to what activities the user is doing before submitting the query. Here focusing on the user's search history which is kept in log format and records the user's past queries and choose the results to view. For example if the user search the queries as "debugging" means the user likes programming. For example if the user searched with the same query "python" before and click the python language on the website link. we have high confidence that the user do it again this time.

Search history have two types short term and longterm.

Short term search history is limited to single search session. It consist sequence of searches with a coherent information need and provides a short period of time. Once a user sends an initial query, check the returned documents if unsatisfied user modifies the query and repeats the search process click through data in the short term search history provide implicit feedback that is used to estimate more accurate query language.

Long term search history unlimited in time scope and all search activities in the past. In long term search history there is no need to detect session boundary.

B. Topic- Sensitive PageRanking Mechanism :

More topics are discussed on the web so , it is difficult to the search engine know what the user want. Personalized search solves this problem in the web search community. First propose a user model to collect the users interest[1]. Users interest can be learned through their click history[7][8]. Next based on this model we have to design an learning method that is used to identify the user interest by analyzing user's click history[2]. To investigate the problem of learning the users interest based on their past click history. Using formal user model develop a method to hide interest automatically based on their past click history. Based on the formal user model develop a method to hidden interest automatically based on past click behavior.

Finally a ranking mechanism that considers a user's hidden interest in ranking pages. The PageRank is based on the random surfer model which essentially models a user doing a random walk on the web PageRank) scheme is an extension of PageRank that provides different types of rankings for different queries[8]. Using Topic Sensitive PageRank instead of Single global PageRank. During the query time the search engine select the most relevant Topic Sensitive PageRank value for the given query and user.

The key idea of PageRank is when a page u links to a page v the owner of page u thinks that page v is important. If a page links to a lot of pages, the important score to each of them are decreased. This simple intuition leads to the formula of computing PageRank. Another explanation of PageRank is based on the random surfer model. The user starts from a random page on the web and the user randomly chooses an out-link to follow.

1. Topic Preference Vector:

A user's topic preference vector means it is a m -tuple $T=[T(1), \dots, T(m)]$ in which m is the number of topics and $T(i)$ denotes the user's degree of interest in the i th topic. The vector T is normalized such that $\sum_{i=1}^m T(i)=1$.

2. Page Preference Vector:

A user's page preference vector means it is a n -tuple $P=[P(1), \dots, P(m)]$ in which n is the total number of web pages and $P(i)$ represents the user's degree of interest in the i th page. The vector P is normalized such that $\sum_{i=1}^n P(i)=1$.

The best way of measured the accuracy is to estimate the user's topic preferences from the real life data collected and ask the users the results are how accurate. The problem of this method is users could tell which topics they are most interested in it is very difficult to assign an accurate weight to each of these topics.

3. Generation of topic preference vector:

In the implementation , if the user is interested in the number of topics it will be fixed to k as an experimental parameter. Then we choose randomly k topics and assign random weights to each selected topic.

4. Generation of click history:

To generate a user's topic preference vector, generate a sequence of L click by the user .

C. Click based and Profile based Personalization Strategies:

Personalization has been proposed for many years it is still problem of consistently effective on different queries for different users. For example A biologist can use the query "mouse" to take information about the biological term, but the programmer can use the same query to find information about computer peripherals. To solve this problem to develop a large scale evaluation framework for personalized search using query logs.[1] Then evaluating five personalized search strategies. That is two click based and profile based ones. These results will be analyzed. Here click based personalization strategies perform consistently and considerably[1] but the profile based ones are unstable. Then



download the search results from MSN search engine. Rerank the results are done by any personalization strategy.

Different kinds of personalized reranking strategies[8] are simulated in this framework. So, the user clicks recorded in query logs automatically. In this framework using click through data recorded in query logs when a user issues a query checks the documents in the result list from top to bottom. The user clicks one or more documents which look more relevant to him and skip the documents which is not interested. The click through data [7][8] will be collected at low cost, so it is possible to do large scale evaluation focusing this framework. In Re-ranking Evaluation Framework then downloading search results from search engine, then the results will be reranked using any one personalization strategy. The click-through data recorded is used to evaluate the re-ranking performance.

The top 50 search results will be downloaded from MSN search engine for the test query. Compute a personalized score for each page using personalization algorithm and then generate a new rank list. Combine the rankings[8] using Borda's ranking fusion method and sort the web pages with the combined rankings. Click based Personalization performs consistent only work on the repeated queries on the same user.

D. Privacy-Enhancing Personalized Web Search:

Personalized search is used to improve the search quality. Here to achieve the balance between users privacy and search quality. Using split and buildup algorithm to collect, summarize and organize their personal information into a hierarchical user profile. Similar terms and parent child terms these are two heuristic rules are used to build a hierarchical user profile in the order of top down fashion. The general terms are ranked to higher levels than specific terms. Private information will be collected using mindetail, expRatio. Mindetail provides which part of user profile is protected and expRatio is used to measure the total amount of privacy with respect to mindetail value.

The specified user profile will be send to the search engine, which is on the server side. Search engine wrapper personalizes the search results with respect to partial user profile. MSN search is chosen as the search engine.

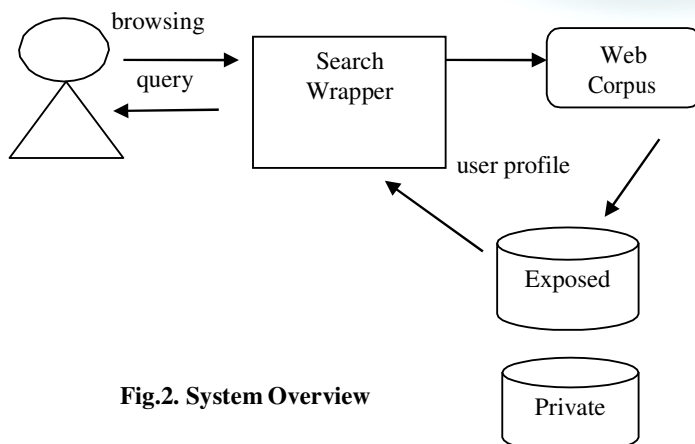


Fig.2. System Overview

That means each result comprises a set of links related to the query and each link gives a rank from MSN search called MSN Rank. User profiles is represented by a weighted term vector. The hierarchical structure is implicitly constructed in a topdown fashion.

[Fig 2]A hierarchical user profile will be automatically created by a user using an algorithm[10]. It represents the user's implicit personal interests and needs. General interests are placed on a higher level, specific interests are placed on a lower level. Only generalized profile will be sent to the search engine with respect to user's own privacy settings. On the server side the search engine wrapper will be developed. Rankings from both partial user profile and search engine results are combined.

Two heuristic rules used in the approach

1. Similar terms:

Two terms are cover the document sets with heavy overlaps can indicate the similar interest.

2. Parent-Child terms:

Specific terms often present together with general terms but the reverse is not true.

In this method deal with unstructured data includes personal documents it is still an open problem on how to define privacy. we improve the balance between privacy protection and search quality.

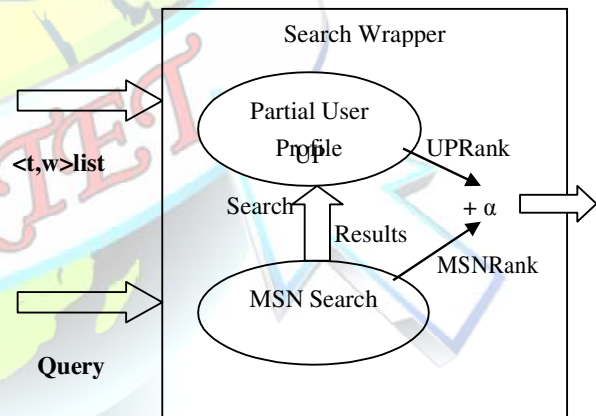


Fig .3. The workflow in the search wrapper

E. Preserving User's Privacy In Web Search Engines:

Web search engines are used to find the huge amount of information in less amount of time. Google is the most used web search engine in the US. Google improve its performance by storing a record of used sites and the user's past searches submitted by each user. Those searches can reveal a lot of information from individual users[1].



For example when a person is applying for a job. At that time if the employer knows that the applicant looking for information for a certain disease, the employee use this knowledge to choose another person for the job. when submitting a query to a web search engine while preserving the user's privacy called private Information Retrieval problem[11].

In these PIR protocol a user can retrieve a certain value from a database which the server holds the database .But it has no knowledge about the data requested by the user. The first PIR protocol scheme is based on several servers kept the same database. These servers unable to communicate between them. The disadvantage of this protocol is it cannot work in single server. Christo Ananth et al. [4] proposed a secure hash message authentication code. A secure hash message authentication code to avoid certificate revocation list checking is proposed for vehicular ad hoc networks (VANETs). The group signature scheme is widely used in VANETs for secure communication, the existing systems based on group signature scheme provides verification delay in certificate revocation list checking. In order to overcome this delay this paper uses a Hash message authentication code (HMAC). It is used to avoid time consuming CRL checking and it also ensures the integrity of messages. The Hash message authentication code and digital signature algorithm are used to make it more secure . In this scheme the group private keys are distributed by the roadside units (RSUs) and it also manages the vehicles in a localized manner. Finally, cooperative message authentication is used among entities, in which each vehicle only needs to verify a small number of messages, thus greatly alleviating the authentication burden.

The client application is a java applet. It is accessed by an html web page that allows users to make a search in a transparent manner. And the client interface displays a form field where the user must type their query. Once the search button is pressed ,the search process will be started. A javascript code is used to sent the text from the form field to the applet. Using this approach the web search engine unable to generate a real profile of individual. It have a moderate delay in the query response. These protocols will not be used in location based services.

Greedy DP algorithm:

Here Greedy DP algorithm used in UPS.It supports online profiling based on predictive metrics of personalization utility and privacy risk.It works in a bottomup manner. The main problem of this algorithm is that it requires recomputation of all candidate profiles which is generated from pruneleaf manner.

GreedyIL algorithm:

The Greedy IL algorithm improves the efficiency using heuristics based on several conclusions.Prune leaf operations reduces the discriminating power of the profile.Greedy Il reduces this measure with heuristic.

III. PROPOSED SYSTEM

Meta Search Engine:

The proposed PWS system is used to extract and learn a user's history and content preferences based on the user's profiles.Different users need the same information and share the same search task.Different users will be valuable for capturing semantically similar search tasks in a global context.To compare task,session and query trails for three search applications determining user satisfaction,predicting user search interests and query suggestion.Sparse query task clustering is used to measure user satisfaction and page utility.Bounded query clustering get information to measure topic similarity for predicting user search interests.

Meta Search Engine is a search tool that uses other searchengine's data and produce their own results from the internet.Search Engine is created and take input from a user and send the queries to third party search engines for results.If data is collected assigned by their ranks and presented to the users.Meta Search Engine is used to enhance the user's experience by consisting multiple results from different search engines.It is used to generating a large volume of data.Spamming also reduce the accuracy of the search.Metasearch engine overcome this problem.There are many types of search engine savy search engine and meta seek engine.The advantage of these engine is sending multiple queries to several other search engines.This is used to extends the search coverage of the topic.Here more search results will be retrieved but using same amount of exertion.It reduces the user's individual type searches.Here searches can be analyzed in three levels a) query level b)quest level c) session level.In query level it fails to capture the interleaving relationships between different quest.A query trail is a sequence of user behavior.It represents all user activities within that particular task URL clicks.Session means a series of queries by a single user made within a small range of time.These three represents determining user satisfaction,predicting user search interest and query suggestion.

There are two types of user behavior search behavior Search behavior and browser behavior. Search behavior is everything the user enters in the search engine to search for the information needed. Browser behavior involves surfing user types a URL address in the browser,king a bookmark or forward page in the browser etc.

Content based recommendation algorithm that estimates the user's interest in unknown content by matching her profile to metadata descriptions of the content.



If a personalized query is identified during generalization the entire runtime profiling will be aborted and the query will be sent to server without a user profile.

Clustering algorithm is used to group documents together based on their similarities thus documents relating to a certain topic will hopefully be placed in a single cluster. This can help users more easily and in getting an overview of the retrieved document set.

Tasktrail is introduced to understand user search behaviors. Task can be defined by a set of semantically relevant queries to satisfy the information need.

A tasktrail represents all user activities within the particular task URL clicks. By separating the different tasktrails from a session it can be used in several search applications such as determine user satisfaction, predicting user search interests and suggesting related queries.

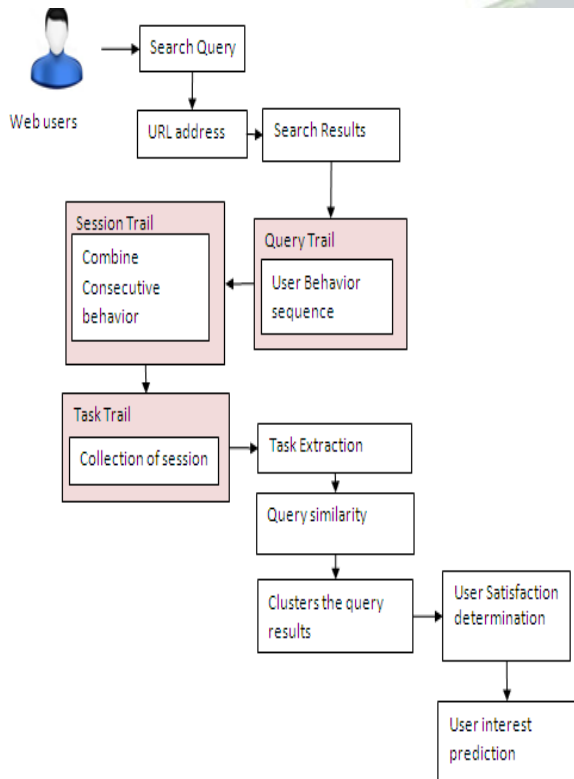


Fig 4. System Architecture

A.Task Definition:

Web log includes users search queries and related actions on the internet. It contains a set of users and it has a sequence of consecutive behaviors e_1, e_2, \dots, e_n where each e_i is a search behavior or a browse behavior. A search behavior is a single query which issued to a search engine. A browser behavior includes one of the following activities 1)user starts

to search from the homepage of the web browser. 2)user types a URL address in the browser. 3)user pastes the URL address from other place into browser. 4)user select a bookmark 5)user clicks the back or forward button in the browser.

B.Task Extraction:

A task is a set of semantically related query trails within a session. If they satisfy these following things 1)they are identical 2) one is a part of the other 3)two partially agree to each other 4) one is a typo of the other the two queries can be grouped into same task. These rules can be used in the annotation process and propose an efficient clustering framework to group queries into tasks. The basic ideas of our clustering framework are described as follows. First, since tasks are extracted out from each session, we follow the time threshold method to segment logs into sessions by choosing a time threshold u . We quantitatively compute the similarity between any two queries. A classifier can be used to learn the weights of various features of query similarity function.

C.Query Similarity:

Classification algorithm is used to compute the similarity between two queries. A labeled data set should be constructed to learn a good query similarity function for task classification. The labels include same task and different task. Features are used to measure the similarity between queries. These features can be categorized into two groups such as time based (temporal) and query word based. Based on query similarity we can construct semantic similarity to extract the results meaning fully. Semantic analysis is nothing but a process of filtering that progressively eliminates more and more input strings until you are left with only valid data.

D.Clustering Queries into Tasks:

From the intuition that consecutive queries more likely belong to the same task than non-consecutive ones, a better approximation is to compute the pair-wise similarity for all consecutive query pairs. In this work, a clustering algorithm Query Clustering using Modified Bounded Spread method (QC-MBSP) is proposed for task extraction.

User Satisfaction Determination:

After the search process, to understand whether a user was satisfied or not in search process, several indirect feedback signals can be used as measures.

- Clicks: The total number of clicks to perform a particular task can be taken as a signal of user satisfaction on that task. Clicking on search results often indicates the relevance between queries and clicked pages.
- Dwell time: Dwell time can also be considered as a signal of user satisfaction. It is because users are more likely to stay on useful pages.



F. User Interest Prediction:

User search interest are represented by queries. queries will be summarized into topics is used to understand the user search interests. Suppose two queries issued by one user they may come from (1) different sessions (inter-sessions); (2) same session (intra-session); (3) different tasks in different sessions (inter-tasks among sessions); (4) different tasks in same session (inter-tasks within sessions); (5) same task in same session (intra-task). All these five sources can provide query pairs. Besides, capturing user search interests at topic level is useful to understand user behaviors

IV CONCLUSION

A major problem in web search is that the interactions between the users and search engines are limited by the small form factors of the search engines and various categories. As a result, web users tend to submit shorter, hence, more ambiguous queries compared to their web search counterparts. We proposed PWS to extract and learn a user's history and content preferences based on the user's profiles. Web log segmentation can be done in query level, session level, or task level. Task trail is nothing but it is a sequence of user behaviors occurred within one session, where they collectively define atomic user information need. Following task trail is an effective method to segment the web log and also to determine the user search behavior. Using time threshold Web logs are segmented into sessions. Queries similar to each other are clustered into same task after computing the query similarity. From the extracted tasks, user search behavior can be determined. The final result display unique web results for unique users based on their searched behavior. We also tested a term matching variant that alleviated the need for an exact term match between queries and trails, which led to coverage and diversity gains at the cost of a slight decrease in relevance.

V REFERENCES

- [1] J. Castelli-Roca, A. Viejo, and J. Herrera-Joancomartí, "Preserving User's Privacy in Web Search Engines," *Computer Comm.*, vol. 32, no. 13/14, pp. 1541-1551, 2009
- [2] B. Tan, X. Shen, and C. Zhai, "Mining Long-Term Search History to Improve Search Accuracy," *Proc. ACM SIGKDD Int'l Conf. Knowledge Discovery and Data Mining (KDD)*, 2006.
- [3] Y. Xu, K. Wang, B. Zhang, and Z. Chen, "Privacy-Enhancing Personalized Web Search," *Proc. 16th Int'l Conf. Wrlld Wide Web (WWW)*, pp. 591-600, 2007.
- [4] Christo Ananth, M. Danya Priyadharshini, "A Secure Hash Message Authentication Code to avoid Certificate Revocation list Checking in Vehicular Adhoc networks", *International Journal of Applied Engineering Research (IJAER)*, Volume 10, Special Issue 2, 2015, (1250-1254)
- [5] F. Qiu and J. Cho, "Automatic Identification of User Interest for Personalized Search," *Proc. 15th Int'l Conf. World Wide Web (WWW)*, pp. 727-736, 2006.

- [6] E. Gabrilovich and S. Markovich, "Overcoming the Brittleness Bottleneck Using Wikipedia: Enhancing Text Categorization with Encyclopedic Knowledge," *Proc. 21st Nat'l Conf. Artificial Intelligence (AAAI)*, 2006.
- [7] K. Ramanathan, J. Giraudi, and A. Gupta, "Creating Hierarchical User Profiles Using Wikipedia," *HP Labs*, 2008.
- [8] K. Järvelin and J. Kekaälä, "IR Evaluation Methods for Retrieving Highly Relevant Documents," *Proc. 23rd Ann. Int'l ACM SIGIR Conf. Research and Development Information Retrieval (SIGIR)*, pp. 41-48, 2000.
- [9] R. Baeza-Yates and B. Ribeiro-Neto, *Modern Information Retrieval*. Addison Wesley Longman, 1999.
- [10] A. Viejo and J. Castellà-Roca, "Using Social Networks to Distort Users' Profiles Generated by Web Search Engines," *Computer Networks*, vol. 54, no. 9, pp. 1343-1357, 2010.
- [11] Lidian Shou, He Bai, Ke Chen, and Gang Chen "Supporting Privacy Protection in Personalized Web Search"