# Web Search Results Personalization: Taxonomy and Current State of Art

**K.Srinivas, Department of CSE, Geethanjali College of Engineering & Technology, Hyderabad, Telangana, India. katkamsrinu@gmail.com**

**V.Valli Kumari, Professor, Department of CS & SE, Andhra University College of Engineering, Andhra University, Visakhapatnam, Andhra Pradesh, India. valli_kumari@rediffmail.com**

**A.Govardhan, Professor, Department of CSE, Jawaharlal Nehru Technological University, Hyderabad, Telangana, India. govardhan_cse@yahoo.co.in**

**Abstract: In the present scenario information can be achieved easily through web. Although search engines provide plenty of data but it has a drawback that the information is not correct. According to research users lack interest if it is time consuming and does not provide accurate data. So search engines must provide accurate data. Since personalized web environments provide exact and instant data both to short term and long term users there is an increase in its demand. We have given various personalization techniques and approaches in this paper.**

**Keywords: Page ranking, Personalization, Search Engine, Web Search.**

## I. INTRODUCTION

Search engines mainly aim at providing accurate data through personalized websites. It also aims at providing browsers to avail the search engine applications like multimedia browsing, visual results and localization immediately on devices like mobiles. According to Anand Spink, a professor at the University of Pittsburgh's school of Sciences and a long time researcher of web search behavior, "Each technology relates to a different aspect of people's information behaviors" which means to response different questions various technologies must be developed. He also mentions in his book 'Web Search': publishing searching of the web (Springer 2004), that there isn't much improvement in search engines from 1997 till now. Since there is not much development, this field has much scope which can be proved by the words "Our challenge is to read a users mind" of Daniel Read's. A problem of limited words is faced not only in search engines but also in localization and personalization of web information; we can overcome this drawback by R&D projects which can provide exact data. For the improvement of Ask Jeeves and Google research is in progress.

### A. Making It Personal

A new attribute of personalization of a search engine will provide exact data needed by the user within less time without displaying unwanted data. As the user lacks facts about search engines he does not know the necessities of it. There is a need to improve search engines so that the users can directly access the data. As per Spink "computer software lacks in knowing users accurate demand and suggests that both implicit and explicit skill must be improved to achieve valuable result."

Two new methods are being introduced by Google so that exact result can be obtained by the user. Depending on the users demand accurate outcome is provided and the user is informed with any new data available through emails using the personalized web search (http://labs.google.com/personalized) which is undergoing beta testing. A site-favored Google search is introduced which is in beta testing, where a particular website is provided with a best suitable report signified by web masters.

Besides merit in question and answer technology, Ask Jeeves is trying hard to be finest in personalization. My Jeeves a new technology provides My Jeeves folder where the browser can store his search outcome and personal data which can be accessed from any PC. Registration is needed for documents above 1000 and if registered there is no limit. In order to provide users with accurate results the organization is scheduling to unite My Jeeves to search skills.

The browser finds the data through hyper textual space, which is uneasy as it provides crowded data which is not accurate to the browsers need. Similar outcomes are provided to the users whose queries are parallel in nature but are presented in different manner. So to overcome this problem personalization can be used which provides accurate results by knowing the needs of the user.

Users mainly use three methods to search data. First is searching data with the use of query and reference given by the user and communicate with other users result with parallel need provides exact outcome. It includes movies, music and products [1, 2]. Another method is browsing through which user can get required data but this is time consuming as the links provide plenty of data which can be irrelevant to the user because of which user needs to go through every link and get the accurate data.

With the use of Information Retrieval (IR) outcome to the query given to the search engine is obtained. A search engine chooses a website which provides appropriate result to the users. The obtained result is used for advance investigation which is dependent on few sources which are not updated.

IR is dependent on order of distinct browser queries causing another difficulty. Information Filtering (IF) is the

only way to overcome the above difficulty depending on the users demand. IF is considered as the best tool as it provides the accurate data depending on the users need and it considers users demand to be constant. So the user is provided with efficient and accurate data. To know whether the data needed by the user is appropriate or not is the biggest task. The disadvantage of IF is that the updating is time consuming because of which it is not much used. We will talk about the main IF modeling tools in the forthcoming papers.

The users can easily obtain data through two methods Information Retrieval where the data is updated and Personalization method is more capable then prior. New tool and its functioning are mentioned below.

Using the above mentioned techniques the user can easily obtain the data through surfing. But the data available is plenty so by using query the required data is obtained. Since the outcome is parallel to all the queries, to overcome this personalization can be used which provides accurate result as per users needs. The data needed by the user is to be mentioned in a sequence so that the data required can be matched with other users whose needs, preferences are same so that accurate result is acquired. Many search engines exists but very few meet users requirements and as competence improves with personalization problems also arise because of which utilization becomes difficult [3] and is a very hard job. The user needs, purpose and interests are stored in external search system and are not linked to many which causes basis for the difficulty. Disadvantage of personalization is that it is time consuming when compared to general and obtaining the accurate results is not that easy. Users are informed with personalization tools, procedures and methods in the upcoming.

## II. OVERVIEW ON PERSONALIZED SEARCH

Above we have discussed about the need and importance of search engines now new tool personalization introduced will be explained. Firstly new slang and terminology will be discussed, then how the user details are used to find user needs from various regions. Lastly different personalized search methods overview is discussed.

### A. Personalization Techniques

Vector Space Model (VSM) [4,5] is dependent on the data existing in IR tool. Every new web page is updated in the web world with a title name, name being from the data mentioned. The outcome for unprepared query is given using these keywords [6]. Search becomes time consuming if the user is not aware of his requirements and does not mention appropriate words, whose reverse is also possible [7]. As per research outcomes user uses 2 to 3 words to find his needs. Search engine faces difficulty with slang also, as a word has synonyms and poli-semis the website may not provide the total information of the synonym used by the user. As a result the search engines using key-word method faces difficulty [8]. Sometimes unconnected data from various web pages is obtained or missed because of synonyms and poli-semis [9] which are time consuming. Since the use of key words has

drawbacks, the search engines must make use of the entire data to provide accurate results to the given query. Principal of content-based personalization approaches is shown in figure 1.

If the need, preferences of different users match then the information is interchanged this is known as Collaborating Method [10, 11]. Social Navigation is the other name proposed by Dieberger et al [12], which is a software where in users can write their views and can vote so that other users can utilize this while searching data by query or surfing.

### B. User Modeling in Personalized systems:

Personalization of search engine uses modeling tool where the information of user and its requirements are stored so that the user is provided with accurate data. During information retrieval or filtering user-modeling method is used. The users are provided with filtered and accurate data as per users needs using this method. Outcomes are good if a personalized search engine is more complex.

This method follows a procedure which includes either registration or simple questionnaire where interests and needs of user are known. But in complex one users need for this data, academic status and his knowledge about the data are known.
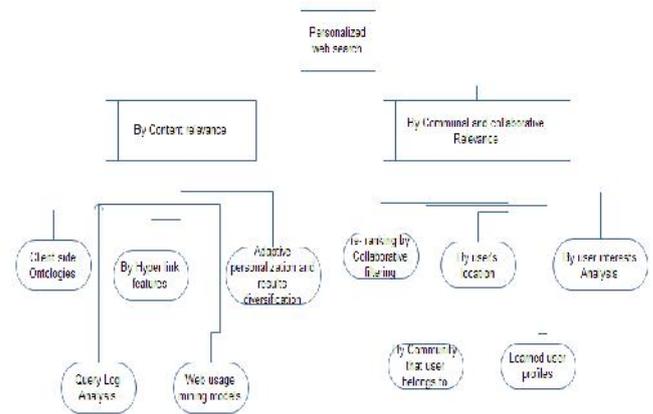


Figure 1: Search Results personalization approaches

## III. INFORMATION DEPENDENCT

This method is dependent on data not the user hence there is no need of knowing users profile and provides filtered and accurate results to the user. Yan named this method as SIFT [13] and experimented this on Internet news articles. SIFT consists of list of words that can and can't be used for searching articles and can be updated by user as per his needs. Every single user is provided with not less than twenty articles and the user can access this data through www browser or if he is unable to connect to web an email is sent to him consisting of the required data.

### A. *Client side ontology based personalization*

In this method outcomes of the user needs (psychology needs) are grouped depending on the concepts. These groups are to be updated and maintained properly. User Conceptual Index [UCI] a new personalized search index was proposed by S. Sendhi Kumara & T.V. Geetha. Personalized Ontology an automatically identified user profile is used in this method and to improve the existing personalized web search engine Page Ontology is used.

### B. *Query log analysis*

Query logs have been analyzed many times [14, 16, 20, 17]. Query distribution depending on query length and query frequency [19, 15, 17], and query type and topicality [15, 20] has been published.

User queries are grouped for the purpose of exploring personalization opportunities. In general personalization has much work as well as in particular personalization to aid web search.

### C. *Personalization which works on hyper link technique*

Weighting links are used at page rank vectors for personalization. URL is used to compute these links, so that particular users profile must match. Binary vectors represent users interests in the form of DNS tree nodes where every feature is connected. On the features of hyperlinks, i.e. anchor terms or URL tokens this method relies. To filter the data user profiles are linked with hyper links.

Depending on the match between the URL and the profile weighted page rank is calculated giving a weight to every URL. Two methods are followed for analyzing web pages.

(i) With the help of graph structure every page is ranked, where a group of research people who concentrate on retrieval tools rank the pages, for example depending on number of hits faced by the page ranking is done, and while results such ranked pages are considered first.

(ii) With the help of hyperlinks URL and anchor texts. The idea of a page is presented by anchor text. The quality and accuracy of the page can be known through URL. The organizer can be known with the use of URL author. To get better content the author must be expert in that particular field, resulting in accurate, good and prompt outcome to the search.

Page Rank method gives likely results for a search and is very useful in personalization. Global ranking to the web documents is done through page ranking, which is done using hyperlinks. Link will be ranked well if it is preferred by many documents. The author of a particular web page is considered as pages with good authors are given importance in page rank. Search engines like Google 1 use web structure for page ranking. To develop accuracy the search results are graded using these ranks. Using page-rank method the value of a web page among billions of documents is estimated. Page ranking is not without lapses [21-23]. The data useful for one may not be useful for another, so before ranking user preferences should also be considered. Sometimes user needs may not be answered well by good page ranks so to personalize page ranks URLs of a page can be utilized nicely.

To overcome the problems of page-rank, a hyperlink feature Internet domains was under taken. User chosen profile is considered as binary vector and is communicated to a DNS tree node and the page rank scores depending on match between hyperlink and profile demands are computed by assigning weights. However the best result is obtained, two problems are faced at the implementation stage. One the user may not be interested in giving his profile and two previous search track are considered to rank the results. Users needs are estimated using the past search track and user profile which may not be appropriate and users do not support this. The user is presented with most appropriate result using ranking method to grade the data.

### D. *Personalization based web usage mining*

Web mining refers estimating and discovering of the data present in www. Web mining is explained in detail in. At present web mining related areas are introduced at later stage WUM is explained

.

Web-mining is linked to the data and below mentioned are its features:

Content data: User is provided with the text and extracting data from web pages is content mining [21].

Structure data: the arrangement of web information system is explained by metadata. The selection of information from the data structure is web structure mining [22,23].

Usage data: User gives the information and then it is connected to the web. As mentioned earlier WUM refers to estimating and discovering users access to the web information system using the present data to modify the web for user was not any new idea but was suggested long back in 1995 [24,25,26]. A good analysis of the research in WUM is given in [27]. The two tiers of WUM Tracking & analysis the present methods are explained here.

### i)*User-interaction tracking*

For personalization the data about the contact of user with internet is used. The data can be obtained in various forms such as web browser on the client side, web server logs or representative server logs. The accuracy of tracking the exact data is as important as an increase in importance of personalization. Data can be stored in web so that the browsers can easily access the page which was used recently. Proxy server logs cannot save total cache hits, effecting the analyzation of search manner and users choice. To solve the trouble an "access pattern collection server" was introduced by Lin et al. (1999), it functions when no secrets are maintained by user. To get the data about hit back stored references Cookey et al. (1999) [28] used referrer and agent fields of server log. Results of various methods were analyzed by Spiliopoulor et al. (2003) [29] and it was observed that server and proxy logs cannot provide user communications temporal aspects. Network-transmitting time is provided to

the logs for document demands where time stamps are stored. The easy access to the data is not possible because of abandoned functioning of the network.

The data is considered to be real and spatial which is present with the user about the contact done with internet. The URL or resource of data can be found out easily because of the total data present with the user, in case of proxy or server logs it's a difficult challenge. Formerly single person used to collect data about the web page for a proxy but now all the users must deliver. It is powerfully done at the users side and is known as session identification. The documents must be separated and must be arranged and set according to their key words for good understanding and examination. To find the user contacts a remote agent was engaged in shahabi et al. (1997) [30]. To collect data from the clients, the run agents are provided with Java scripts or Java applets and the user may not like to integrate Java program in the browser. An illustration about data collecting process depending on the user-side collecting idea was given by Shahabi et al.(2000) [31].

## ii) Access pattern analysis

Since the data is large its investigation is difficult. During access models construction few features cannot be considered. The number of hits paper has faced by the user projects its value or grade [32]. Amongst the outcome of document browsing first document is chosen. To know this characteristic and future position the non independent aggregate tree and hidden Markov models are employed. In the milieu of web personalization applications temporal features like page view time are given importance besides spatial features. As per Yan et al.(1996) [32] and lovene and logic (2000) the judgment of paper cannot be done depending on its selection time because it may happen that few papers are not selected due to its difficult accessing method, Zipfin division. To overcome this problem the view time is combined with other features. The combination of many features is explained in the present model. As per the requirement of situation either complexity of the model or accuracy can be sacrificed because of availability of various features.

To find out the features various data-mining methods and approaches were offered by scientists. According to Mobasher et al [33] in personalization of web grouping methods provide best outcome to that of group regulations. Depending on the previous browser selections to picture future situation other non independent methods are used. The important similarities among page selections are understood and signified by these methods and a Markov method is introduced for this purpose by Zukerman et al [18] and Cadez et al. Yan et al. [34] commenced the grouping to mine usage data method. Each part of the vector represents importance of a characteristic, for instance hit-count for correlating to the web page in the vector structure. To get user access technique group algorithm is used.

To access the grouping achievement in the milieu of WUM many algorithms were tested. Cluster miner a new grouping algorithm was introduced to respond particular web

personalization necessities by Perkowitz and Etzioni. A capable hierarchical algorithm was proposed by Fu et al. and BIRCH. On the belief of fuzzy data usage Joshi and Krishnapuram [37] prefer a fuzzy relational clustering algorithm for WUM. The correlation of cluster mining achievement with two other grouping methods which are vibrant in machine learning research was given by Paliouras et al [39]. A single cluster is divided depending on the users features collected by different groups were spotted by Mobasher et al [33]. User's dynamic profiles in connection with static profiles anonymous WUM was tested by Verder Meer et al [40]. To prepare group model for updating new developments in user's behavior, Dynamic clusters method is explained.

## iii) Search results usage mining

To attain web usage mining for search results personalization two steps process method is used.

### a) Tracking

Firstly the spatial and temporal characters of user interactions are traced and remote ages are used to attain it.

### b) Analysis

The use of features that are taken as medium and are assigned to the available information are used for estimation. At different stages the quality of periods are the qualities of its parts, information is gathered at different stages of each part. With the use different parts the estimation of periods is calculated. Similarity is evaluated to prove the 'similarity' as the last step.

## E. Adaptive personalization

This mainly focuses on personalization for the purpose of user necessities and is primarily for the user. Personalization involves three steps.

- The viewed data is found.
- The current demands are provided exactly.
- The temporary data needed is supplied. It has been verified that this process has the capability of doing the above three steps.

Personalizing a web search can be done in various ways and each has its respective advantage. One of the ways is knowing users view on the outcome which the users express indirectly.

This model is flexible, depends on above mentioned three steps and is not preset on any characteristic. Personalization depends on importance given to the data and the search behavior in the past. The importance is given to users for accurate data. The importance is either given to the general users ranking or to the current situational want. This method includes both the characteristics and consequently they are used.

The other procedure of personalization conveys importance and searches results both together. The search outcomes are presented by finding similarities between them. Though there are changes in the user's needs and page ranking, this method is considered to be the best. As per the users conduct and needs this process is improved. As per the

topic when users are to be related personalization prefers more than one browser.

### i)Refining known information

To review the outcome again back button is used indicating users demand for that data. Users choose to rewrite the query again for the same outcome rather than referring it using back button. If the query is viewed in the past and wants the same in the present it indicates user wishes to review the previous outcome.

There are three ways to find results for such desire. (i) To get the previous outcome the user appeals the local index to suggest the exact query. The search agent modifies the query and presents to the search engine making a developed query. The user got the control with him and can change the query at various stages. In this model the outcome of the last query is considered to be the best. (ii) As per the viewed pages data the query is transformed accordingly. Pages with best data are not much preferred, the main reasons can be either the page displaying users keywords is not appropriate or the user possesses idea of the page so does not prefer it. (iii) The user is linked with other page through a click the result can be obtained.

### ii)Finding out about topics of user interest

The user's interest can be known if he is a regular user or is a long way searcher, in this case if the user inputs incorrect query search engine can help the user by providing the accurate data by tracking his old searches. The user gets good result when the same topic is viewed again instead of new topic.

Depending on the users interests the outcomes are provided where in if a query is raised for the first time or if it does not have any data. The keywords or abstract are used to structure the query. Present search outcome includes the search outcome of around 10 pages of old search outcome, providing the user with previous search records results.

### iii)Serving an Ad-Hoc information need

The user may desire many topics but few of them may for temporary purpose. The above model considers the first query only and as per the search results the query is updated.
The user needs and preferences are improved with the keywords of the selected outcome, helping in categorizing the given results. The query term near to the selected data is improved by adding key word from the selected data. At present the user can establish new query using the similar old query words. If the same query is asked or if the user wants to view the previous query then an illustration of the previous queries in the present is very important.

### F.Limits of the models discussed under content relation based approach

i)Data restrictions: IR tools are limited to certain characteristics of data like text and image.
ii)Over-specialization: since the data provided is as per users demand so the data is filtered and accurate to users' needs.

iii)Communal/collaborative: Besides personalizing user needs and past search field context is also used for searching.

In two ways searching can be done depending on more than one basis. (i) The users whose demands are same and has similar context are grouped through search engines special characteristic. (ii) The query data can be referred easily as we know that for query qt page pj will be selected from a group of selected users. The main advantage of collaborative search is it does not depend on past analysis instead it focuses on personalizing search results of users group. No rules are fixed for search engines in collaborative search as it can use the contents with web pages, graphics and photos, audios and video. Ranking metric is used to examine the number of times the page is used.

## IV. BY COMMUNAL AND COLLABORATIVE RELEVANCE

In this section we discuss the approaches like collaborative re ranking, group based also known as user's community, user's location based and user interests based approaches.

### A. Collaborative re-ranking of search results

In this method the user and society preferences are combined to provide accurate outcomes, as a result the data is limited to their preferences. It helps in providing exact data and then grouping. Collaborative method is used to rank the outcomes and assigning weights to the user and society profiles by providing inputs to the search and meta search engines.

### i)Features of this architecture

The combined preferences of the user and society are expanded by the characteristics of architecture.
Characteristics are as mentioned below:
- User and society data collections are associated and thus pages are classified.
- Meta search engines are offered.
- Browsers data is stored.
- Society is presented.
- Clusters are maintained.
- To search data pre processing is applied.
- To rank the combined set of outcomes post processing is used.

Besides this other facilities are offered so that the user can obtain accurate outcome to the query presented. For personalization all or few characteristics are used and are illustrated below.

### ii)Preprocessing for context search

Context can be used to obtain exact outcome from plenty of data. Context refers to users' needs, importance and the nature of community. Through questionnaire users interest can be known where he selects the answer which perfectly suits him indicating his needs.

### iii)Processes of profiling users

Weight vector for a single user profile is constructed using the step word, the words which best suits user is opted. Since

plenty of data is present to overcome this user must submit the selected data but the user may not be willing to do so and if he does then also he may not provide the entire data. Fixed process is used to calculate word-weight vector. Vector includes group of words { r t } with weights assigned u i w. the words can be sequenced and profiled properly if the user ranks and clicks the data of word-weight vector. User in a group can be considered in user profile, profile here includes user profile and the data opted by him is also updated. If the page is recommended then the problem of finding groups with similar features will not occur.

### iv)Process of profiling group

Word-weight vector is contained in profiling group; from the selected groups these words are opted. As per the position of users the group for word-weight vector is decided.

### v)Processes of managing community

The aims, needs and preferences are matched with the selected data in a community. One problem in profiling is that there can be no matches. Second problem is to collect the data about the companies having similar user needs groups. Similar groups and data that suit the groups are to discovered by the administrator. The administrator will be overloaded if the groups do not take any role. On the groups and selected data profile is reliable, so it has to be prepared properly. It is an advantage as the method works though group profile is not present. It becomes difficult to deal with users when groups are not present.

### vi)After search process (post processing for context search)

Outcome of the search data are ranked by post processor and is done in two ways. (i) The matter or pointer of the data is matched with the previous data present with the group and the data is ranked well if there is any match. (ii) With the use of word-weight vectors the appropriate profile is asked based on the content. Data can be assessed if the key word of the data and profile match.

### vii)Process of updating user profile weight

If the user is obtaining accurate outcome to its query then the user or group is updated with that data which helps in organizing and ranking user profiles. Set consists of the words that are most frequently used by the user. Using Rochio tool that depends on Rel list the relatedness of a search query can be assessed. Standard Rochio tool is used other than relevance feedback system with Rel list. The use of non-relevant words is the main difference among the present and old model. The relevant words are not considered by the present model since the data does not have certain or standard resource. The group profile and ranks of the experts are improved automatically by a transition of user profiles.

### B.　Group based approach

### i)Repetition and regularity in search groups

Searching for accurate outcomes is usually done individually. The tools are considered as a common task and the user search indirectly. For instance if an author is keen in knowing the tools and equipments of work place, he is not only provided with the data but also with the search boxes at the entry itself. To gain standard web searches the boxes will

have the data that have been viewed by other users as per their interests. The user is connected to the portal form and a temporary group who have common needs in work place equipments which is an important characteristic to all the search engines. The outcomes of the portal are more related to the given topic and the other things being the same. The groups can be formed through search boxes if the interests of different people match and come in contact with those of the writer, sales man or a school research. Since these groups have same search interests and behavior they are more useful. For instance, fig 2 illustrates the result of an H-week research of the search models for a group of about to employees at a local software company. The search queries of more than 20,000 individuals and nearly 16,000 search results were estimated.
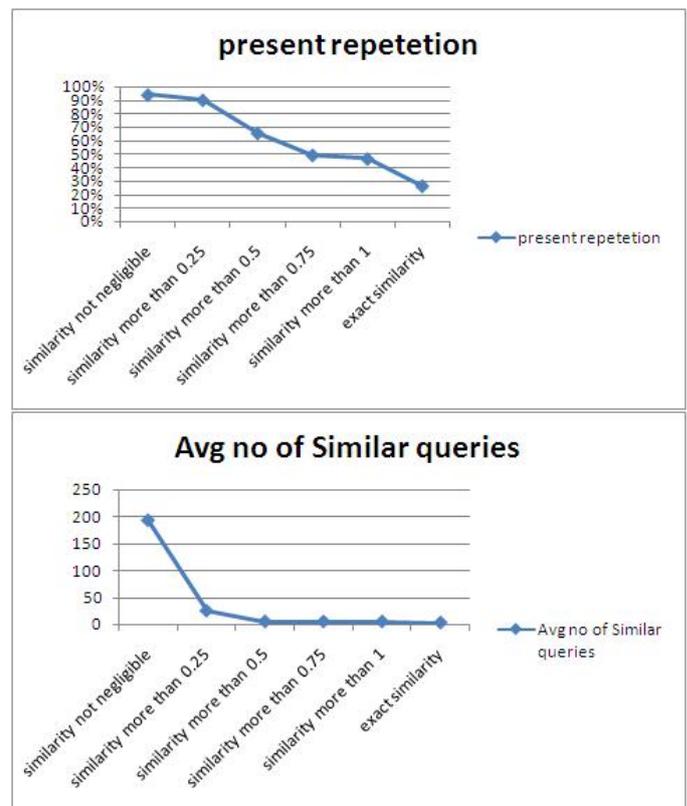


Figure 2.Search methods.

The above figure describes the Statistics of a 17-week research of the group models for a group of about 70 employees at a local software company and the percentage of the repetition of query with parallel words and similar queries.

The average similarity between the queries during the research is observed in the Figure 2. On average, about 65 percent of given queries had 50 percent of same query words (0.5 similarity threshold) along with five other queries at maximum; above 90 percent of given queries maximum 25% of similar words were shared with  25 other queries at minimum. So the users find the temporary corporate search group in the similar manner when compared to generic search scenarios resulting in less repetition rates of about 10 percent at the 0.5 similarity threshold.

The results of similar other research search society, society of users have same needs and preferences support that the outcome of web search is a repetitive and regular activity. Depending on the users query and outcome selected a type of community search knowledge is formed, suggesting the possible way to attach the sharing of users searching experiences between the society group and members. For instance, if a visitor to wildlife portal is searching for "jaguar pictures," the search engine can suggest search results of the similar queries which were asked by the other members of the society, which helps to know the similarities among the members of the society for wildlife. This helps the trainee searchers to gain the shared knowledge from the experienced searchers of the society.

The research results as shown above depict the percentage of similarity of interests and needs of the members of the society which are resistant by other research results. The result of such research tells us that the sharing of search data among the members of the society must be educated to the society. So instead of wasting time on long search processing the search results can be obtained from the other members of the group who has already viewed the same query.

## ii)Collaborative web search

By storing the past results of the user queries the past or underlying knowledge of the queries and their results at the community level can be shared by the search communities. The main objective here is to protect replacing of predictable search engine by providing the relevant information to the society without default outcomes. For instance the queries related to "jaguar", the communities must be provided with the wild life related pages to that of the related query. Thus the relevant outcomes from the society are related to the top of the page results where as other results are simply related but are not positioned at the top.

## a) Capturing community search knowledge

The behavior of the members of the society towards their query selection, outcome and their frequent usage is stored and is termed as capturing a community's search behavior. Populating the community search matrix HC called a hit-matrix, such that HC ij represents the number of times the result page pj has been viewed for the similar selected query qi is conceptualized. Each row of community's hit-matrix refers to the number of times the result is selected for the searches that have been made for several times by the members of the group for the same qi where as the column of the community's hit-matrix refers to the number of times the community members have selected the same pj for different qi.

## b) Making relevant promotions

How it is possible that the present query, qT, will be used from a group's hit-matrix as potential promotion candidates to find the results? With respect to the qT the group history data must be identified the past qT. The hit-matrix assumes that such pages are present and are most frequently selected with respect to the qT and the data can be used to know the importance of such pages. For instance, the importance of

result pj with respect to the qT and the relative proportion of selections that the pj has got for this qT and is shown in equation 1 below:

$$\text{Re}levance^C(p_j, q_T) = \frac{H_{T_j}^C}{\sum_{\forall j} H_{T_j}^C} \quad (1)$$

The query-relevance methods limit the users who are selected for the improvement of the pages that were previously selected for the specific target query qT. As mentioned in fig 2 only 25 percent of query requests in the test society matches to that of past demand. For instance, a direct way to calculate query similarity is by computing the proportion of terms shared by qT and some other query qi and is shown in equation 2 below:

$$sim(q_T, q_i) = \left| \frac{q_T \cap q_i}{q_T \cup q_i} \right| \quad (2)$$

As shown in equation 3 query-similarity metric can be used as the ground for a modified relevance metric:

$$W \ \text{Re}l^C (p_i \ q_T \ q_1 \ ...q_n) = $$
$$\frac{\sum_{i=1..n} \text{Re}levance^C(p_i, q_i).sim(q_T, q_i)}{\sum_{i=1..n} Exists^C(p_i, q_i).sim(q_T, q_i)} \quad (3)$$

The relation between the pj and the qT is calculated by independently computing the exact query importance of page pj with respect to a set of queries (q1, ..., qn) which are assumed to be similar to qT. It is observed that the queries that share 50 percent of its terms with the target query are only measured. The importance of pj with respect to qT is weighted to the sum of the individual exact query importance values, as the relation between the pj and the qi is discounted by the similarity of $q_i$ to $q_T$. As a result the most frequently used pages for query similar to qT are preferred when compared to the less frequently used pages that are selected for query similar to qT.

## C. Location based approach

## i) Personalization by query rewriting

By expanding the original query by users location to personalize search results is a sensitive way. For instance the users query is search "Chinese restaurants" and user's location is "newyork", then a new query can be framed like "Chinese restaurants Newyork" by relating the location to the query the user can get accurate result, the new query is submitted to the search engine so that the user can be provided with exact outcome. The search engine will provide the user with Chinese restaurants located in newyork as the new query consists the location which is helping the user to find appropriate result. Because of two errors this method has failed: (1) The aim of the query may not only be the users location or because of limited accuracy of the location the result may not exist.(2)The location of the user may be determined imperfectly.

## ii)Personalization by re-ranking

Page re-ranking is done to weight users' location in a traditional method. Re-ranking is done by selecting the top K data with original query and after the ranks of the data are improved the search results are reorganized by matching the users' location. As a result the data that match users' location are ranked higher and are kept at the top than those data that do not match the users' location. In the top data the user locations match are separated and then weighted for various segments. The scores of the data after re-ranking sr(q; d; l) for query q, Web document d and user location l is as shown below:

$$s^r(q,d,l) = s(q,d) + \sum_i w_i I_i(d,l) \quad (5)$$

here s(q; d) represents the original rank score i.e. before personalization, Ii(d; l) represents a pointer function that shows if user location l exist in data segment i and wi is the weighting parameter for data segment i. To maximize the importance after personalization supervised learning is used to assess the parameter wi.

### D. Search results personalization based on user profiles

By utilizing user profiles that have various abstraction stages i.e. from general to specific stage which can present various environments for personalization re-ranking was analyzed. A system that allows the automatic improvement of structured user profiles, which are developed based on a available category hierarchy was proposed by Pretschner and Gauch [43]. A user profile as a weighted concept hierarchy, which is developed from the Open Directory Project (ODP) was suggested by Operetta and Gauch [44]. ODP was used to recognize user profiles for personalized web search by Siege et al. [45]. Since 590,000 models/ideas are available in ODP only a few models are used that are at the top level in the ODP hierarchy and the models at the lower level in ODP hierarchy are ignored. This may affect the ranking of the individuals whose features are not presented in the high-level categories in ODP. Since all high-level categories are present in user profiles, the user profile may contain many unrelated categories on using prevailing hierarchy, so to beat these disadvantages of utilizing an existing taxonomy/hierarchy, Kim and Chan [46] originated a way to raise a user interest hierarchy (UIH) by understanding from implied user behavior. Soon using a scoring function [47, 48] for personalized ranking with the UIH this was further improved. Based on user profile and the outcomes obtained from search engines a page can be ranked. Web pages in users' bookmarks and the achieving task mostly based on DHC is used to enlarge UIH users profile.

### E. Analysis of user interests and activities

To personalize web search many efforts were made in the past. Personalization can be done by knowing the general concern of the users. For instance Google asked the users to prepare their profile by grouping their needs and concern. To personalize search results by mapping web pages to the same group by utilizing users profile. Gauch et al. [50, 51] has discovered personalize web search outcomes and this method was used by commercial information-filtering systems. To construct a personalized version of page rank [53] for locating the past query-independent on web pages, personal profiles were utilized in the web search background. Parallel tool was utilized for mapping user queries to groups based on previous search history outcomes by Liu et al. [53].

Using the tools like relevant feedback or query enhancement the data of the users' goals can be gathered. Various interface tools were tested for varying transparency to know from the users how query can be expanded by Koenemann and Belk [54]. For producing query refinements different alternative approaches were explained by Antic [55] and McKeon et al. [56]. A very short-term approach of a users concern was attached with relevant data and query refinement and it is involved that first the query must be submitted and then it must be modified accordingly. Explicit query refinement is seldom used in general especially in Web search [55]. Users are least interested in specifying their concerns in general. As per the recent survey by Tee van et al. [56] he recommended that users are not interested in specifying their needs and concern instead they directly search for their demands by browsing via pages or by specifying their needs in the form of queries. Nielsen [57] feels the same to that of the above result and says that for personalization users need not take extra pain. No appropriate outcomes are possible when people are forced to take extra effort on indicating their needs and concern [58].

Jaime Tee van et al [59] introduced an automated analysis of activities and needs of users for search results personalization. To personalize users' current web search, search procedures are prepared in which users past relations with a wide variety of content is considered. This tool emphasizes on trying to gather data about users need and concern instead of depending on the postulation that users will mention about their needs and is utilized to re-rank Web search results within a relevant feedback structure.

## V. CONCLUSION:

To overcome the data overload problem personalized search on the Web is a research field and it is being concerned by the people. The cause over here is that data plays an important part for the users as the users face the problem of getting outcomes to their search results which is important to the user to achieve both personal and professional tasks. The problems like searching problem, search engines accuracy and the time consumed by the users to get accurate outcome for the given query can be reduced by the talent to filter and create a personalized set of resources.
The originality and spirit of the personalization field imply that, in future new and exciting procedures and methods will be projected and will be provided to the users so that the users can relate to his daily use like search engines or desk top search tools. The two essential research fields anthology and semantic web have started being noticed in this framework.

### REFERENCES

[1] Resnick, P., Varian, H.R.: Recommender systems. Commun. ACM 40(3) (1997) 56–58
[2] Montaner, M., Lopez, B., Rosa, and J.L.D.L.: A taxonomy of recommender agents on the internet. Artificial Intelligence Review 19 (2003) 285–330

[3] Khopkar, Y., Spink, A., Giles, C.L., Shah, P., Debnath, S.: Search engine personalization: An exploratory study. First Monday 8(7) (2003) http://www.firstmonday.org/issues/issue8_7/khopkar/index.html.

[4] Salton, G., McGill, M.: An Introduction to modern information retrieval. Mc-Graw-Hill, New York, NY (1983)

[5] Rijsbergen, C.J.V.: Information Retrieval. Butterworth-Heinemann, Newton, MA,USA (1979)

[6] Micarelli, A., Sciarrone, F., Marinilli, M.: Web document modeling. In Brusilovsky, P., Kobsa, A., Nejdl, W., eds.: The Adaptive Web: Methods and Strategies of Web Personalization. Volume 4321 of Lecture Notes in Computer Science. Springer-Verlag, Berlin Heidelberg New York (2007) this volume

[7] Olston, C., Chi, E.H.: ScentTrails: Integrating browsing and searching on the web. ACM Transactions on Computer-Human Interaction 10(3) (2003) 177–197

[8] Furnas, G.W., Landauer, T.K., Gomez, L.M., Dumais, S.T.: The vocabulary problem in human-system communication. Commun. ACM 30(11) (1987) 964–971

[9] Freyne, J., Smyth, B.: An experiment in social search. In Bra, P.D., Nejdl, W., eds.: Adaptive Hypermedia and AdaptiveWeb-Based Systems, Third International Conference, AH 2004, Eindhoven, The Netherlands, August 23-26, 2004, Proceedings. Volume 3137 of Lecture Notes in Computer Science., Springer (2004) 95–103

[10] Goldberg, D., Nichols, D., Oki, B.M., Terry, D.: Using collaborative filtering to weave an information tapestry. Commun. ACM 35(12) (1992) 61–70

[11] Resnick, P., Iacovou, N., Suchak, M., Bergstrom, P., Riedl, J.: Grouplens: an open architecture for collaborative filtering of netnews. In: CSCW '94: Proceedings of the 1994 ACM conference on Computer supported cooperative work, New York, NY, USA, ACM Press (1994) 175–186

[12] Dieberger, A., Dourish, P., H¨o¨ok, K., Resnick, P., Wexelblat, A.: Social navigation: techniques for building more usable systems. Interactions 7(6) (2000) 36–45

[13] Kritikopoulos, A., Sideri, M.: The compass filter: Search engine result personalization using web communities. In Mobasher, B., Anand, S.S., eds.: Intelligent Techniques forWeb Personalization, IJCAI 2003Workshop, ITWP 2003, Acapulco, Mexico, August 11, 2003, Revised Selected Papers. Volume 3169 of Lecture Notes in Computer Science., Springer (2003) 229–240

[14] S. M. Beitzel, E. Jensen, A. Chowhury, D. Grossman, and O. Frieder. Hourly analysis of a very large topically categorized web query log. In SIGIR, 2004.

[15] A. Broder. Taxonomy of web search. SIGIR Forum, 36(2), 2002.

[16] B. Jansen and U. Pooch. A review of web searching studies and a framework for future research. Journal of American Society for Information Science and Technology, 52(3), 2001.

[17] B. Jansen and A. Spink. How are we searching the web? a comparison of nine search engine query logs. Information Processing and Management, 42, 2006.

[18] T. Kuflik and P. Shoval. Generation of user profiles for information filtering - research agenda (poster). In SIGIR, 2000.

[19] C. Silverstein, M. Henzinger, H. Marais, and M. Moricz. Analysis of a very large search engine query log. SIGIR Forum, 33(1), 1999.

[20] A. Spink and B. Jansen. A study of web search trends. Webology, 1(2), 2004.

[21] Cohen, W., A. McCallum, D. Quass. 2000. Learning to understand the web. IEEE Data Engrg. Bull. 23 17–24.

[22] Kuo, Y. H., M. H. Wong. 2000. Web document classification based on hyperlinks and document semantics. PRICAI 2000 Workshop on Text and Web Mining, Melbourne, Australia, 44–51.

[23] Henzinger, M. 2000. Link analysis in web information retrieval. Bull. of the Technical Committee on Data Engrg., IEEE Computer Soc. 23 3–9.

[24] Armstrong, R., D. Freitag, T. Joachims, T. Mitchell. 1995. Web-Watcher: A learning apprentice for the world wide web. AAAI Spring Sympos. on Inform. Gathering from Heterogeneous, Distributed Environments, Stanford, CA, 6–13.

[25] Lieberman, H. 1995. Letizia: An agent that assists web browsing. Proc. of the Internat. Joint Conf. on Artificial Intelligence, Montreal, Canada, 924–929.

[26] Pazzani, M., L. Nguyen, S. Mantik. 1995. Learning from hotlists and coldlists: Towards a WWW information filtering and seeking agent. Proc. of IEEE Internat. Conf. on Tools with AI, Washington, DC, 39–46.

[27] Srivastava, J., R. Cooley, M. Deshpande, P. N. Tan. 2000. Web usage mining: Discovery and applications of usage patterns from web data. SIGKDD Explorations 1 12–23.

[28] Cooley, R., B. Mobasher, J. Srivastava. 1999. Data preparation for mining world wide web browsing patterns. Knowledge and Inform. Systems 1 5–32.

[29] Spiliopoulou, M., B. Mobasher, B. Berendt, M. Nakagawa. 2003. Evaluating the quality of data preparation heuristics in web usage analysis. INFORMS J. on Comput. 15(2) 171–190.

[30] Shahabi, C., A. M. Zarkesh, J. Adibi, V. Shah. 1997. Knowledge discovery from users web page navigation. Proc. of the IEEE RIDE97 Workshop, Birmingham, England, U.K., 20–31.

[31] Shahabi, C., F. Banaei-Kashani, J. Faruque. 2001. A reliable, efficient, and scalable system for web usage data acquisition. WebKDD'01 Workshop, ACM-SIGKDD 2001, San Francisco, CA. http://dimlab.usc.edu/Research.html.

[32] Yan, T. W., M. Jacobsen, H. Garcia-Molina, U. Dayal. 1996. from user access patterns to dynamic hypertext linking. Fifth Internet. World Wide Web Conf., Paris, France, 1007–1118.

[33] Mobasher, B., H. Dai, T. Luo, M. Nakagawa, Y. Sun, J. Wiltshire. 2000a. Discovery of aggregate usage profiles for web personalization. Proc. of the Web Mining for E-Commerce Workshop WebKDD'2000, Boston, MA. http://maya.cs.depaul.edu/mobasher/papers/webkdd2000/webkdd2000.html.

[34] Breese, J. S., D. Heckerman, C. Kadie. 1998. Empirical analysis of predictive algorithms for collaborative filtering. Proc. of ncertainty in Artificial Intelligence, Madison, WI, Morgan Kaufmann, San Francisco, CA.

[35] Fu, Y., K. Sandhu, M. Shih. 1999. Clustering of web users based on access patterns. International Workshop on Web Usage Analysis and User Profiling (WEBKDD'99), San Diego, CA, 18–25.

[36] Zhang, T., R. Ramakrishnan, M. Livny. 1996. BIRCH: An efficient data clustering method for very large databases. SIGMOD '96, Montreal, Canada, 103–114.

[37] Joshi, A., R. Krishnapuram. 1999. Robust fuzzy clustering methods to support web mining. Proc. of SIGMOD Workshop in Data Mining and Knowledge Discovery, Seattle, WA, 1–8.

[38] Strehl, A., J. Ghosh. 2003. Relationship-based clustering and visualization for high dimensional data mining. INFORMS J. on Comput. 15(2) 208–230.

[39] Paliouras, G., C. Papatheodorou, V. Karkaletsis, C.D. Spyropoulos. 2000. Clustering the users of large web sites into communities. Proc. of the Internat. Conf. on Machine Learning, Stanford, CA, 719–726.

[40] VanderMeer, D., K. Dutta, A. Datta, K. Ramamritham, S. B. Navanthe. 2000. Enabling scalable online personalization on the web. Proc. of the 2nd ACM Conf. on Electronic Commerce, Minneapolis, MN, 185–196.

[41] X. Shen, B. Tan, and C. Zhai. Implicit user modeling for personalized search. In CIKM, 2005.

[42] K. Sugiyama, K. Hatano, and M. Yoshikawa. Adaptive web search based on user profile constructed without any effort from users. In WWW, 2004.

[43] A. Pretschner and S. Gauch. Ontology based personalized search. In Proc. 11th Intl Conf. on Tools with AI, pages 391–398, 1999.

[44] M. Speretta and S. Gauch. Personalized search based on user search histories. In Proc. Intl. Conf. Web Intelligence, pages 622–628, 2005.

[45] A. Sieg, B. Mobasher, and R. Burke. A large-scale evaluation and analysis of personalized search strategies. In Proc. CIKM, pages 525–534, 2007.

[46] H. Kim and P. Chan. Learning implicit user interest hierarchy for context in personalization. In Proc. Intl. Conf. on Intelligent User Interfaces, pages 101–108, 2003.

[47] H. Kim and P. Chan. Personalized ranking of search results with learned user interest hierarchies from bookmarks. In O. Nasraoui, O. Zaine, M. Spiliopolou, B. Mobasher, B. Masand, and P. Yu, editors, Web Mining and Web Usage Analysis, pages 158–176. Springer, 2006.

[48] Personalized Web Search by Using Learned User Profiles in Re-ranking, Jia Hu and Philip K. Chan, http://www.cs.fit.edu/~pkc/

[49] Google Personal. http://labs.google.com/personalized

[50] Gauch, S., Chafee, J. and Pretschner, A. (2004). Ontologybased personalized search and browsing. Web Intelligence and Agent Systems, 1(3-4): 219-234.

[51] Speretta, M. and Gauch, S. (2004). Personalizing search based on user search history. Submitted to CIKM '04. http://www.ittc.ku.edu/keyconcept/

[52] Jeh, G. and Widom, J. (2003). Scaling personalized Web search. In Proceedings of WWW '03, 271-279.

[53] Liu, F., Yu, C. and Meng, W. (2002). Personalized Web search by mapping user queries to categories. In

[54] Proceedings of CIKM '02, 558-565.

[55] Koenmann, J. and Belkin, N. (1996). A case for interaction: A study of interactive information retrieval behavior and effectiveness. In Proceedings of CHI '96, 205-212.

[56] Anick, P. (2004). Using terminological feedback for Web search refinement: a log-based study. In Proceedings of WWW '04, 89-95.

[57] McKeown, K. R., Elhadad, N. and Hatzivassiloglou, V. (2003). Leveraging a common representation for personalized search and summarization in a medical digital library. In Proceedings of ICDL '03, 159-170.

[58] Teevan, J., Alvarado, C., Ackerman, M. S. and Karger, D. R. (2004). The perfect search engine is not enough: A study of orienteering behavior in directed search. In Proceedings of CHI '04, 415-422.

[59] Nielsen, J. Personalization is overrated. In Jakob Nielsen's Alertbox for October 4, 1998. http://www.useit.com/alertbox/981004.html.

[60] Personalizing Search via Automated Analysis of Interests and Activities, Jaime Teevan MIT, CSAIL, Susan T. Dumais Microsoft Research, Eric Horvitz Microsoft Research