# Personalized Web Page Recommendation System with Diversified Ranking

## Priyanka M[1], Nidhi B[2], Namita S[3] , Karthikayani[4]

[1]Student, Computer Science, SRM University, Chennai , Tamilnadu
[2]Student, Computer Science, SRM University, Chennai, Tamilnadu
[3]Student, Computer Science, SRM University, Chennai, Tamilnadu
[4]Asst.Professor, Computer Science, SRM University, Chennai, Tamilnadu

--------------------------------------------------------------------------***---------------------------------------------------------------------------

**Abstract -** *As of today, the web has become the largest inventory of knowledge which allows users to access the most relevant information within no time, which then makes it paramount to organize data in accordance with users needs. A well organized web service recommendation involves satisfying the users potential requirements based on the usage history. In this paper we propose a web service which employs a very effective and precise algorithm to cluster relevant and related queries based on association of mined words. After which we adopt a ranking algorithm which provides diversity and excludes redundant information to portray only admissible data. Finally, we compute the scores based on the personal usage history and potential users interest.*

***Key Words*:**  Cluster, Relevant information, Ranking, Diversity, Mined Words

## 1. INTRODUCTION

Search engines are basically programs that use keywords to search for documents that relate to these keywords and then puts the words in order of relevance to the topic that was searched for. They are used to filter the information that is on the internet and transform it into results that each individual can easily access and use it within a matter of seconds.



**Figure 1** : Personalized user profile search along with browsing history

The primary goals of a modern day search engine include effectiveness and efficiency. Effectiveness implies on the quality of the search results,i.e to retrieve the most relevant set of documents for a query. We can process text and store text statistics to improve search results.

The latter, efficiency focuses on the speed at which information retrieval occurs as it is paramount to process the queries from the users as fast as possible. To perform this we use specialised data structures. Other goals include scalability, adaptability and incorporation of new and relevant data.

Clustering of web search results is an attempt to organise the results into a number of thematic groups in the manner a web directory does it. This approach, however, differs from the human-made directories in many aspects. First of all, only documents that match the query are considered while building the topical groups. Clustering is thus preformed after the documents matching the query are identified. Consequently, the set of thematic categories is not fixed – they are created dynamically depending on the actual documents found in the results. Secondly, as the clustering interface is part of a search engine, the assignment of documents to groups must be done efficiently and on-line. For this reason it is difficult to download the full text of each document from the Web. Clustering ought to be performed based solely on the snippets returned by the search service.

In recent years, Web clustering engines have been proposed as a solution to the issue of lessening lexical ambiguity in Web Information Retrieval. These systems group search results, by providing a cluster for each specific aspect (i.e., meaning) of the input query. Users can then select the cluster(s) and the pages therein that contain the best answer their query needs. However, many Web clustering engines group search results on the basis of their lexical similarity, and therefore suffer from synonymy (same query expressed with different words) and polysemy (different user needs expressed with the same word).
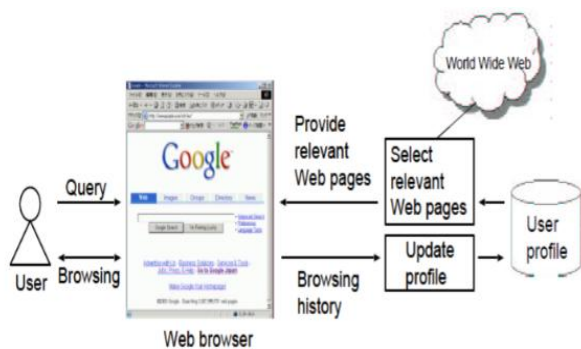
## 2. RELATED WORK:

Web service recommendation is a process of discovering and recommending suitable Web services to  users. Various works have been done on service recommendation based on

quality of service (QoS). Most of them employed Collaborative Filtering (CF) techniques, some of them applied content-based approach, and a few of them combined CF approach with content-based techniques. Quality-of-Service (QoS) is widely employed to represent the non-functional characteristics of Web services and has been considered as the key factor in service selection. Quality of service (QoS) is the overall performance of a computer network, particular performance seen by users. Different user's obtain completely different QoS values when invoking the same Web service, and  the observed QoS values are smaller or larger than the corresponding value released by the provider of the Web service. Therefore, due to the underlying assumption that service consumers tend to obtain the best recommendations from those with similar QoS preferences or usage experiences to themselves, personalized QoS-aware Web service recommendation appears as an emerging technique to address the above issue.

CF (Collaborative Filtering), also known as social filtering or social information filtering, is the most popular mode in the field of personalised recommender systems. It's  aim is to predict and identify the data  (e.g., website, commodity, social networking service, etc.) the  user might be interested in  according  to  historical  data,  and  to  make recommendations on this basis. To  our knowledge, there are basically two types of  CF methods for Web service recommendation:  memory-based  and  model-based approaches. Even with  some criticism, the CF-based methods have been mainly  used in prior studies and in many commercial systems, and their feasibility and good performance have also been validated in terms of different data sets.

In general, the memory-based approaches can be divided into three categories: user-based , item-based , and hybrid approaches . The basic idea of these  types  of approaches is to conduct rating predictions based on historical QoS records after finding out similar users or items. Even if  they are easy to  implement and are cost effective, there are  various drawbacks with these type of approaches, such as the bias of human ratings and the relatively poor scalability with large-scale  data  sets.The next  approach to  clear  the  query ambiguity is search result clustering. Given a query, a flat list of text snippets returned from commonly-available search engine is clustered using some notion of textual similarity. Search result clustering approach can be classified as data-centric or description-centric.

Diversification is another research topic dealing with the issue of query ambiguity. Its aim is to reorder the top search results using a criteria that maximise their diversity. Similarity functions have been used to measure the diversity among the documents and between document and query. Other techniques include the use of conditional probabilities to determine which document is most different from higher-ranking ones or use affinity ranking, based on topic variance and coverage.
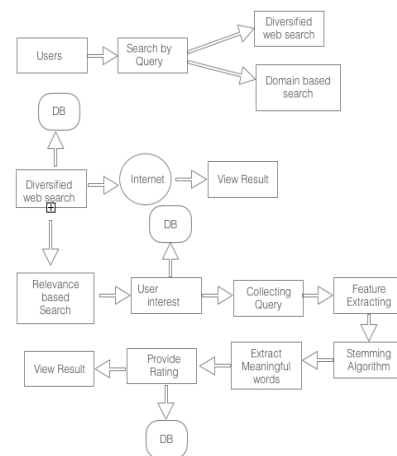
## 3. ARCHITECTURE:



**Figure 2**: Architecture of search engine

Now we describe the architecture of the personalized web service recommendation system which takes diversity into consideration. To provide a personalized search results we should maintain a personal user profile. A user profile stores information about the user interests and preferences. Once the user profile is generated then search is carried forward using either of the two options provided : Diversified web search (or) Recommended search.

Diversified web search   evaluates based on relevance and text similarity, It is used for indexing a search query. The process of diversified web search takes a query and stores the mined words in the database and at the same time it crawls for information in the world wide web. In case of the recommended search the evaluation occurs on the basis of the relevance of the user's historical interest with web services based on a content based similarity measure. The users historical interest can be mined his/her own service usage or query history. The ranking is performed by scoring the user interest and then displaying the top-k results. To collaborate the query results we use clustering of relevant and related url with regard to the association rules. Our approach is based on two criteria: one is the queries themselves, and the other on user clicks. In this framework we basically extract the query from the usage history and collect similar queries and form clusters. Once we cluster the data they are evaluated on the basis of number of clicks in feature extraction process. The feature extraction will analyze the frequency factor and accordingly rank the information.

## 3.1 Query clustering:

User's queries can be classified into different query clusters. Concept based user profiles are employed in the clustering process to obtain the personalization effect. The most similar pair of concept nodes and then, merge the most similar pair of query nodes, and so on. Each individual query submitted by each users treated as an individual node and each query

with a user identifier. We perform the grouping in a similar dynamic manner, whereby we first place the current query and clicks into a query group. A web query topic classification/categorisation is a problem in information science. The task is to assign a Web search query to one or more predefined categories, based on its topics. The importance of query classification is underscored by many services provided by Web search. A direct application is to provide better search result pages for users with interests of different categories. For example, the users issuing a Web query "apple" might expect to see Web pages related to the fruit apple, or they may prefer to see products or news related to the computer company. Online advertisement services can rely on the query classification results to promote different products more accurately. Search result pages can be grouped according to the categories predicted by a query classification algorithm. However, the computation of query classification is non-trivial. Different from the document classification tasks, queries submitted by Web search users are usually short and ambiguous; also the meanings of the queries are evolving over time. Therefore, query topic classification is much more difficult than traditional document classification tasks.

**Table 1:** Query clustering with examples

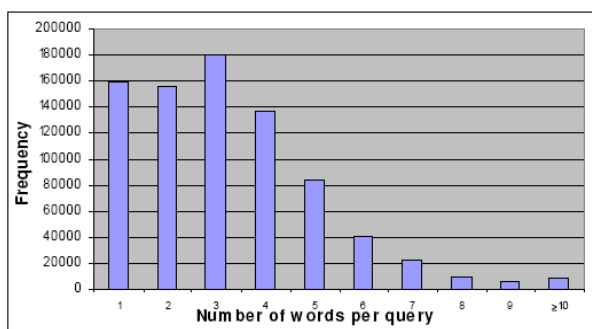| Query | Categories |
|---|---|
| Apple | Computers\Hardware\Living\ Food&Cooking |
| FIFA 2015 | Sports/ Soccer/ schedules & Tickets / Entertainment |
| Cheesecake Recipies | Living \ Food & Cooking \ Arts & Humanities |
| Poem | Humanities \ Living \ Dating & Relationship |



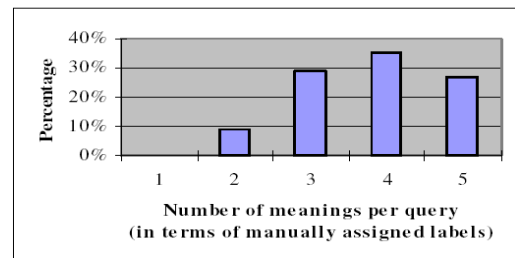**Figure 3**: A graph to show the number of words per query against the frequency



**Figure 4** : A graph to show the number of meanings per query against percentage

## 3.2 Web-Based search:

Web-search queries are typically influenced by several metrics: Content relevance derived from documents' anchor text , title and headings, word frequency and proximity , file, directory, and domain names, and other more sophisticated forms of content analysis, User behaviour extrapolated from user's time-spent- on-page, time-on-domain, clickthrough rates, etc, Popularity in the global link structure with authority , readability , and novelty typically determining the linkage. Links to the most "relevant" pages, according to the above criteria, are then potentially clustered and delivered to users who in turn browse the results to find the desired information. Although researched in detail along most of the mentioned criteria, search engines still leave a lot to be de-sired. In this paper, we emphasize one important inefficiency of state-of-the-art search engines: content redundancy, and propose a system that significantly improves search results for learning-type queries. Looking from a user's perspective, we review an existing classification of Web-search queries that aim at predominantly textual content.

• Navigational – the user is seeking a Web-site with an unknown URL; typically, a single specific Web-site is the "correct answer" to the posed query (e.g., Washington mutual points to www.wamu.com).

•Journey – the user is browsing a certain content category on the Web at random trying to discover further points of interest (e.g., indoor plants).

• Shopping – the user is looking for the best offer on the Web for a product/service (e.g., nikon d300 price).

• Learning – the user wants to know a specific detail or the entire breadth of knowledge available on the Web for a specific query (e.g., rhododendron).

While queries from the first three categories could be handled with relatively simple URL lists: a single "correct"

URL, a randomized list of "relevant" pages, and an unbiased comparative shopper, learning-type queries remain difficult to address. The main culprit is the cumulative information redundancy over the ranked list of returned documents. It dominates search results to the point where finding content that is not displayed on an encyclopedia style Web-page such as Wikipedia, usually ranked at the very top of results, is a cumbersome task that requires browsing dozens of links. In this paper, our goal is to propose a simple yet powerful tool for pass-filtering a small set of text documents that offers the greatest joint coverage on a given topic.

If we denote as SQ the total knowledge that exists on the Web about a given query Q, we want to build a search engine that returns a set of essential pages1 EQ that maximizes the information covered over SQ. While SQ is truly a semantic digest of all Web content related to Q, we argue that a simple "bag-of-words" approach to representing SQ is a surprisingly efficient model. Then, we formalize the overall optimization objective using a weighted coverage function that takes into account both word and page relevance. Using the Sequential Forward Floating Selection (SFFS) algorithm, we show that a fast URL ordering by their joint knowledge coverage is achievable and well accepted by users. Figure below illustrates an abstract example of how SQ is covered by a set of pages computed using a traditional page ranking algorithm (top) and a set of essential pages assembled to maximize their joint query coverage (bottom). As a result, in the traditional model, in order to learn details about SQ users have to browse substantially more pages.

In a simple implementation, we used an existing search engine to obtain a large list of most relevant URLs for a given query, then used our tool for post-processing, i.e., re-ordering, of these results. This way, we avoided building and running a Web crawler and a Web index. In an over-whelming majority of cases, our ranking was substantially different than the ranking returned by the underlying search engine. In a user study with over 120 search queries and 35 subjects, in approximately 5 out of 6 learning queries users found our ranking to provide better or equal learning experience to Google's. This result is impressive from the perspective that we did not use to any other metric from the {C,U,P} set to improve our rankings. Our optimization strategy is orthogonal to the {C,U,P} methods in its objective and can be combined with them in an arbitrary manner to achieve superior page ranking results.
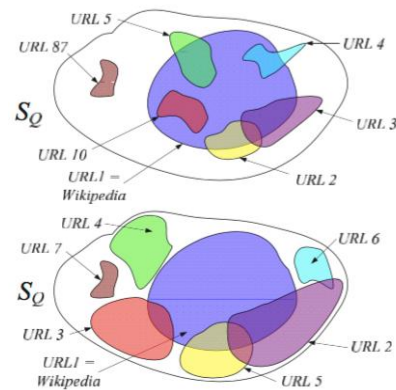


**Figure 5** : An abstract example of the query coverage by an existing (top) and the proposed algorithm(Bottom).

### 3.3 Relevance-Based search :

We first review a traditional relevance-based search mechanism related to our technology. Consider a database D of Nd documents. The goal of relevance-based rank-ordered search is to generate a permutation $\pi Q(D)$ based on a search query Q so that the documents that have higher relevance to Q, come higher in the retrieval results.
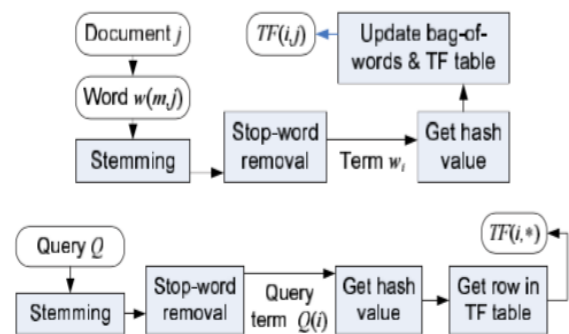


**Figure 6**: Document pre-processing and indexing to generate the term-frequency table (top). Query processing steps in relevance-based search engines (bottom).

Figure 6 shows the rank-ordered search process present in many search engines. For each document j in D ($1 \leq j \leq Nd$), all words in the document are first extracted. The m-th word in the j-th document, w(m,j), then undergoes stemming . In this step, the word root is retained while word endings are removed. Words such as as, is, be, etc., in a pre-defined set of stop-words are then removed as they do not describe the context semantics. Stemming and stopping improve search

performance by giving users more pertinent results; they also reduce the search complexity by reducing the dictionary of words. We denote the total number of unique terms in the resulting list T as Nt. Term frequency TF(i,j) indicates the number of times the i-th term appears in the j-th document. The term frequency information for D and T is organised as a term frequency table of size Nt × Nd. To facilitate fast access, a hash table is constructed to map each term to the corresponding row of the term frequency table.

The term frequency data provides valuable information about document relevance, and is employed as a core variable to define the relevance score in rank-ordering documents. One example metric is the relevance score CW (i, j) , which is defined as:

$$CW(i,j) = \frac{(K_1 + 1)\, CFW(i)\, TF(i,j)}{K_1\, [1 - b + b \cdot NDL(j)] + TF(i,j)}.$$

Here, CFW(i) denotes the cumulative frequency of the i-th term in the entire database, and is obtained as:

$$CFW(i) = -\log_2(n_i/N_d),$$

where ni is the total number of documents in D which contains the i-th term. NDL(j) in Eq. 1 represents the normalised length of the j-th document and is computed by dividing the length of the j-th document, L(j), by the average document length Lavg, i.e., NDL(j) = L(j)/Lavg. Constants K1 and b could be tailored to meet the needs of specific types of databases; some example values are K1 = 2 and b = 0.75. The term frequency table and an adopted relevance score metric are the two main tools used in relevance- based search engines.

A single-word search query Q is handled as follows. First, the query word Q undergoes stemming and stopping. Then, its hash value is used to point to the corresponding row of the term frequency table. Documents in D are then sorted in decreasing order of their relevance score computed using Eq. 1. This process is shown in Figure 2. If a query Q contains multiple terms {Q(1),..., Q(Nq )}, we compute the relevance score R(j) for document j as follows:

$$\mathcal{R}(j) = \sum_{i=1}^{i\,\tau_q} CW(k_i, j).$$

The relevance score R(j) is computed for all documents in D and used to rank order D in response to multi-term query.

## 4. DIVERSYFYING SEARCH RESULTS:

Existing literature has also considered diversity of Web-search results as an additional factor for ordering documents. A re-ranking technique was proposed in [1] based on the maximum marginal relevance criterion to reduce redundancy from search results as well as to present document summarisations. Zhai et al. have defined the subtopic retrieval problem as finding documents that cover as many different subtopics of a general topic as possible [2]. In [3], the authors propose an affinity ranking scheme to re-rank search results by optimising diversity and information richness of the topic and query results.

The aforementioned three techniques [1, 2, 3] model the variance of topics in groups of documents. They all have difficulties applying the concept of diversity to Web-search, e.g., in , the authors assume that all pages are labeled with the topics they cover, then rank them to roughly improve the number of topics covered by a set of pages. Clearly, in the most applicable case of Web-search, such labeling is not available a priori. For that reason, we do not compare our search engine to experimentally, rather we chose to compare our results to state-of-the-art search engines such as Google's, for which we speculate that they address the problem of result diversity. In contrast to prior known methods that focus on maximizing diversity, the technology introduced in this work aims at modeling the overall finite knowledge space for a specific query and improving the coverage of this space by a set of documents. We propose a "sack-of-words" model for representing knowledge spaces, introduce a formal notion of coverage over the "sack-of-words," and derive a simple but systematic algorithm to select documents that maximize coverage, while being relevant to the search topic.

Document indexing for essential pages is equivalent to the process illustrated in Figure 6. For each page on the Web, its distinct set of words is extracted, stemmed using , and filtered for stop words . The global term-frequency table is computed and stored. Given a single-term query Q, the subset of documents, DQ, contain- ing Q is identified using the global term-frequency table. As DQ contains all the information about Q, we denote the set of terms extracted from DQ as SQ.

We informally formulate the problem of essential page selection as finding a subset of documents EQ ⊂ DQ that provides maximum coverage about the query, where the formal notion of information coverage is introduced later. Let NdQ = ||DQ|| and NtQ = ||SQ||. We remark that for a

single- term query, all documents in DQ contain the query term; for queries containing multiple terms {Q(1),…, Q(m)}, at least one of these terms appears in each document in DQ. We denote the subset of the global term-frequency table that relates to the search query Q as TFQ ≡ SQ ×DQ and record its size as NtQ × NdQ. For each term, t ∈ SQ, relevant to the query, we define a term-relevance score, r(t):

$$r(t) = \frac{n_t^Q}{N_d^Q},$$

where nQt represents the number of documents in DQ which contain t. The term-relevance score heuristically measures how relevant t is to Q; the higher the score, the higher the relevance. We use term-relevance as a relevance metric with important notes: it is different from more complex, widely used relevance metrics such as in Eq. 1, it is used due to its simplicity and demonstrated semantic effective- ness while used to identify essential pages.

We define a coverage score, C(j), of a document j ∈ DQ:

$$\mathcal{C}(j) = \sum_{\forall t_i \in \mathbb{S}_Q} \gamma(t_i) \times TF^Q(i, j).$$

where TFQ(i,j) represents the term-frequency value of the i-th term, ti, in document j; γ(ti) ,it quantifies the overall importance given to covering ti in EQ, henceforth we refer to this metric as the term-importance score, and define it as follows:

$$\gamma(t_i) = r(t_i) \log_2 \left[ \frac{1}{r(t_i)} \right].$$

Figure 7 shows the variation of γ(ti) vs. r(ti). The rationale behind choosing such a metric to describe word-importance is as follows:

• Low r(ti) – words that are less relevant to the query do not provide significant information about the query, and therefore they are less important.

• High r(ti) – words that are very relevant to the query (such as the query words itself) provide more information about the query. However, they appear in most documents containing the query word. Hence, it is of less importance to cover these words among the bag- of-words, and therefore they are assigned a low word- importance score. This explains the reasoning behind γ(ti) → 0 as r(ti) → 1.

• Important words – the remaining words are deemed relatively important; our algorithm aims at covering as many as possible of these words with a fixed-cardinality subset of pages from DQ. Intuitively, this is the information that occurs relatively often in DQ and thus is likely to be semantically related to Q. It also occurs relatively infrequently so that it could be spread over a number of pages and thus it can be inconvenient for a user to search for it using a traditional search engine.
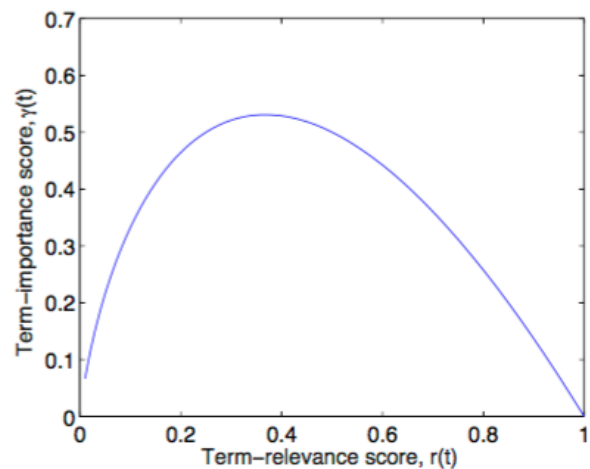


**Figure 7**:Term-importance score as a function of the term-relevance score.

## 5. SURVEY :

In this section, we aim to quantify the observed benefits of the proposed technology. For our survey, we consider several search queries such as ainu people, alexander mackenzie, bill gates, Britney spears, and Mariana trench. Table 2 shows the top 10 most frequent terms in SQ for queries bill gates and Britney spears, respectively. We observe that these stemmed terms are highly relevant to the search query.

**Table 2:List of most frequent stemmed   terms for two exemplary search query**

| Query | bill gates | | britney spears | |
|---|---|---|---|---|
| Order | Dictionary | Freq. | Dictionary | Freq. |
| 1 | gate | 3046 | britnei | 3265 |
| 2 | bill | 2271 | spear | 2806 |
| 3 | microsoft | 787 | music | 464 |
| 4 | busi | 354 | time | 385 |
| 5 | peopl | 325 | video | 382 |
| 6 | comput | 310 | celebr | 304 |
| 7 | softwar | 246 | album | 263 |
| 8 | job | 241 | search | 254 |
| 9 | search | 235 | page | 252 |
| 10 | page | 234 | pop | 246 |

**Table 3: List of most "important" stemmed terms for two**

| Query | bill gates | | britney spears | |
|---|---|---|---|---|
| Order | Dictionary | Freq. | Dictionary | Freq. |
| 1 | live | 90 | singer | 142 |
| 2 | system | 142 | live | 135 |
| 3 | product | 109 | song | 149 |
| 4 | window | 193 | web | 175 |
| 5 | technolog | 128 | custodi | 204 |
| 6 | servic | 65 | peopl | 206 |
| 7 | william | 205 | mtv | 127 |
| 8 | melinda | 141 | boi | 94 |
| 9 | internet | 108 | award | 152 |
| 10 | foundat | 195 | blackout | 165 |

**exemplary search queries.**

The corresponding stemmed terms that have the highest word- importance scores, $\gamma(w)$, as defined are shown in Table 3 for the two queries. For the query bill gates, we notice that words such as live, window , Melinda that are highly relevant to the query and are not covered by all documents, have a high $\gamma(w)$. Similar inferences can be made for the query britney spears from Table 3, suggesting that the proposed word-relevance metric can identify important words.
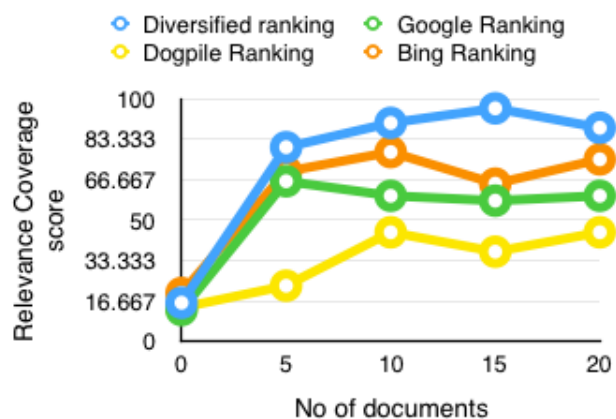


**Figure 8:** Comparison of search ranking algorithms

## 6. CONCLUSION:

In this paper, we describe a fully equipped search engine which will associate relevant data based on user interest with a help of a personalized account. We compare this with previous well established search engines to demonstrate the limitations present in the other search engines, that include:

- Web content indexing does not contain proper structure.

- Ambiguity of information resulting poor correlation.

- Redundancy of information leading to irrelevant data.

- Personalized account, to provide potential user interest.

We have presented a novel query comparing descriptor for measuring the similarity of queries, that incorporates both content and click-through information. This is based on a key insight – search engine results might themselves be used to identify query similarity, and may thus best combine the complementary strategies of content-ignorant and content-aware clustering.

## REFRENCES:

[1] S. Brin and L. Page. The anatomy of a large-scale hypertextual web search engine. Proc. of the WWW, pp.107–117, 1998.

[2] S.E. Robertson and K. Sparck-Jones. Relevance weighting of search terms. Document Retrieval Systems, Vol.3, pp.143–60, 1988.

[3] S. Lawrence and L. Giles. Context and page analysis for improved Web search. IEEE Internet Computing, vol.2, no.4, pp.38–46, 1998.

[4] Ashwin Swaminathan, Cherian Varkey Mathew, and Darko Kirovski. Essential pages. Technical Report MSR-TR-2008-015 January 2008.

[5] Ji-Rong Wen, Zhicheng Dou, Ruihus Song, "Personalized Web Search", Microsoft Research Asia, Beijing, China, 2009.

[6]W. Ng, L. Deng, and D.L. Lee, "Mining User Preference Using Spy Voting for Search Engine Personalization," ACM Tra , Internet Technology, vol. 7, no, 4 article 19,2007.

[7] Ben franklin, Genealogical data Mining,2006 .

[8]D. Beeferman and A. L. Berger, "Agglomerative clustering of a search engine query log," in KDD, pp. 407–416, 2000.

[9]J.-R. Wen, J.-Y. Nie, and H. Zhang, "Query clustering using user logs," ACM Trans. Inf. Syst., vol. 20, no. 1, pp. 59–81, 2002.