# Privacy Protection using Various Algorithms in Personalized Web Search

## Nivi.A.N[1], Vanitha.S[2], Saranya.K.R[3] and Yamini.S[4]

[1, 2, 3, 4] Department of Computer Science and Engineering, Coimbatore, India
annivi1992@gmail.com
vanithaarvind@gmail.com
saranyaravi251291@gmail.com
yamini.shree18@gmail.com

## ABSTRACT

The importance of using a personalised web search will give totally different search results for various users organize search results otherwise for every user, based upon their interests, preferences, and knowledge needs. Personalised web search differs from generic web search that returns same analysis results to all users for identical queries, in spite of varied user interests and knowledge. So, the personalization could be a general would like in net search now-a-days. This paper includes the review of assorted algorithms towards personalization [1].

**Keywords:** Privacy, Personalized web search, Page ranking, Information retrieval

## 1. INTRODUCTION

In fast development of new technologies, search engines plays polar role in data retrieval. Personalized web search can offer totally different search results looking on the user's preference[8]. Personalized web search has completely different levels of effectiveness for various queries, users and context, so one personalization algorithm cannot improve accuracy of ranking for all queries and it may even have an effect on the accuracy of search beneath sure circumstance A good personalization algorithm depends on made user profiles and internet corpus[6]. However, as the web corpus is on the server, re-ranking on the shopper facet is bandwidth intensive as a result of it needs an outsized range of search results transmitted to the shopper before re-ranking. The spectacular growth within the measure of information on the web has attracted an enormous selection of users towards it. Search engines gift a well-organized thanks to search the relevant info from the web.

However, the search results non inheritable won't always be useful to the users, as computer program fails to recognize the user intention behind the question. With the exponential growth of the net, we have a tendency to even have a marketplace for personalized internet search systems. For instance, a keyword search for ‖pen‖ is ambiguous. The user may well be yearning for the piece implement with a degree from that ink flows or a pen — feminine swan[7]. Personal information, i.e. browsing history, emails, etc., are mostly unstructured, that it's onerous to live privacy. Additionally, it is additionally tough to include unstructured information with search engines while not summarization. So, for the aim of each net personalization and privacy preservation, it's necessary for associate degree algorithm to gather, summarize, and organize a user's personal information into a structured user profile[5]. Meanwhile, the notion of privacy is very subjective and depends on the people involved. Things thought-about to be non-public by one person

may be something that others would like to share. during this regard, the user should have management over that components of the user profile is shared with the server. This paper targets at bridging the conflict desires of personalization and privacy protection, and provides varied algorithmic program wherever users decide their own privacy settings supported a structured user profile[9].


# 2. DESCRIPTION OF OUR ALGORITHMS

## 2.1. Page Rank Algorithm

Page Rank is a numeric value that represents how important a page is on the web. Google figures that when one page links to another page, it is effectively casting a vote for the other page. The more votes that are cast for a page, the more important the page must be. Also, the importance of the page that is casting the vote determines how important the vote itself is. Google calculates a page's importance from the votes cast for it. Page Rank is Google's way of deciding a page's importance. It matters because it is one of the factors that determine a page's ranking in the search results. It isn't the only factor that Google uses to rank pages, but it is an important one. A page "votes" an amount of Page Rank onto each page that it links to. The amount of Page Rank that it has to vote with is a little less than its own Page Rank value (its own value * 0.85). This value is shared equally between all the pages that it links to. From this, we could conclude that a link from a page with PR4 and 5 outbound links is worth more than a link from a page with PR8 and 100 outbound links. The Page Rank of a page that links to yours is important but the number of links on that page is also important. The more links there are on a page, the less Page Rank value your page will receive from it. If the Page Rank value differences between PR1, PR2,.....PR10 were equal then that conclusion would hold up, but many people believe that the values between PR1 and PR10 (the maximum) are set on a logarithmic scale, and there is very good reason for believing it. Nobody outside Google knows for sure one way or the other, but the chances are high that the scale is logarithmic, or similar. If so, it means that it takes a lot more additional Page Rank for a

page to move up to the next Page Rank level that it did to move up from the previous Page Rank level. The result is that it reverses the previous conclusion, so that a link from a PR8 page that has lots of outbound links is worth more than a link from a PR4 page that has only a few outbound links. Suppose we have 2 pages, A and B, which link to each other, and neither have any other links of any kind. This is what happens:-

Step 1: Calculate page A's Page Rank from the value of its inbound links

Page A now has a new Page Rank value. The calculation used the value of the inbound link from page B. But page B has an inbound link (from page A) and its new Page Rank value hasn't been worked out yet, so page A's new Page Rank value is based on inaccurate data and can't be accurate.

Step 2: Calculate page B's Page Rank from the value of its inbound links

Page B now has a new Page Rank value, but it can't be accurate because the calculation used the new Page Rank value of the inbound link from page A, which is inaccurate. It's a Catch 22 situation. We can't work out A's Page Rank until we know B's Page Rank, and we can't work out B's Page Rank until we know A's Page Rank. Now that both pages have newly calculated Page Rank values, can't we just run the calculations again to arrive at accurate values? No. We can run the calculations again using the new values and the results will be more accurate, but we will always be using inaccurate values for the calculations, so the results will always be inaccurate. The problem is overcome by repeating the calculations many times. Each time produces slightly more accurate values. In fact, total accuracy can never be achieved because the calculations are always based on inaccurate values. 40 to 50 iterations are sufficient to reach a point where any further iterations wouldn't produce enough of a change to the values to matter. This is precisely what Google does at each update, and it's the reason why the updates take so long. One thing to bear in mind is that the results we get from the calculations are proportions. The figures must then be set against a scale (known only to Google) to arrive at each page's

actual Page Rank. Even so, we can use the calculations to channel the Page Rank within a site around its pages so that certain pages receive a higher proportion of it than others.

## 2.2. Brute Force Algorithm

Brute force could be a simple approach to resolution a retardant, usually directly supported the problem's statement and definitions of the ideas involved. usually it concerned iterating through all potential solutions until a sound one is found. Although it's going to sound imbecilic, in several cases brute force is that the best thanks to go, as we can consider the computer's speed to resolve the matter for North American nation. Brute force algorithms conjointly present a pleasant baseline for North American nation to check our additional complicated algorithms to. Brute-force algorithms don't seem to be typically clever or particularly economical, however they're value considering for many reasons:

• The approach applies to a good type of issues.

• Some brute-force algorithms area unit quite sensible in follow.

• it's going to be a lot of bother than it's value to style and implement a a lot of clever or economical algorithmic program over employing a simple brute-force approach

As an easy example, contemplate ransacking through a sorted list of things for a few target. Brute force would merely begin at the primary item, see if it's the target, and if not consecutive move to succeeding till we tend to either realize the target or hit the tip of the list. for little lists this is often no drawback (and would truly be the well-liked solution), except for extraordinarily giant lists we tend to may use a lot of economical techniques.

## 2.3. The GreedyDP Algorithm

 The first greedy formula GreedyDP works in an exceedingly bottom up way .The unvarying technique terminates once the

profile is generalized to a root-topic. The best-profile-so-far square measure the final results of the rule. the most disadvantage of GreedyDP is that it desires recompilation of all candidate profiles (together with their discriminating power and privacy risk) generated from makes an attempt of prune-leaf. This causes important memory wants and machine worth. GreedyDP formula works on 2 key ingredients they are:

### 2.3.1. Optimal sub-structure

A best answer to the complete downside contains inside it best solutions to sub issues (this is additionally true of dynamic programming)

### 2.3.2. Greedy choice property

 Greedy choice + Optimal sub-structure gives the correctness of the greedy algorithm

## 2.4. The GreedyIL Algorithm

The GreedyIL rule improves the efficiency of the generalization practice heuristics supported several findings. One necessary finding is that any prune-leaf operation reduces the discriminating power of the profile. In alternative words, the exile displays monotonicity by prune-leaf. GreedyIL algorithmic rule selection Properties are:

### 2.4.1. Locally optimal choice

   – Make optimistic choice available at a given moment.

### 2.4.2. Locally optimal choice     globally optimal solution

   – In other words, the selection of greedy is always safe.
   –To prove this algorithm Exchange Argument are used.

### 2.4.3. Contrast with dynamic programming

   – Choice at a given step may be depend on solutions to sub problems (bottom-up)

# 3. COMPARISON

| S.no | Algorithm | Advantages | Disadvantages |
|---|---|---|---|
| 1 | Page Ranking | Page Rank is entirely general and apply to any graph or network in any domain. Page Rank is currently often employed in bibliometrics, social and knowledge network survey, and for link forecast and recommendation. It's even used for analytic thinking of road networks. | It favors the older pages, as a result of a brand new page, even an awfully sensible one won't have several links unless it's an area of AN existing web site. It may be simply inflated by the construct of "link farms". However, whereas looking, the search actively tries to search out the issues. |
| 2 | Brute Force | It yields cheap algorithms for a few necessary issues like sorting, matrix operation, closest-pair, convex-hull. It yields commonplace algorithms for straightforward process tasks and graph traversal issues | It seldom yields economical algorithms. Some brute force algorithms intolerably slow e.g., the algorithmic formula for computing Fibonacci numbers. it's not as constructive or artistic as another style techniques |
| 3 | GreedyDP | It works in a very bottom-up manner. The iterative technique terminates once the profile is generalized to a root-topic. The best-profile-so-far square measure the ultimate results of the rule. | It wants re computation of all candidate profiles (together with their discriminating power and privacy risk) generated from makes an attempt of prune-leaf. This causes vital memory wants and machine value. |
| 4 | GreedyIL | It improves the efficiency of the generalization practice heuristics supported several findings. One necessary finding is that any prune-leaf operation reduces the discriminating power of the profile | The GreedyIL rule alternative properties at a given step could depend upon solutions to sub issues (bottom-up). |

Table-1 Comparison of Privacy Protection Algorithms

## 4. CHALLENGES OF PERSONALIZED SEARCH

Despite the attractiveness of personalized search, there is no large-scale use of personalized search services currently. Personalized net search faces many challenges that retard its real-world large-scale applications:

1. Privacy is a problem. Personalized net search, especially server-side implement, needs grouping and aggregating plenty of user info together with query and click through history. A user profile can reveal an oversized quantity of personal user info, such as hobbies, vocation, financial gain level, and political inclination that is clearly a heavy concern for users [4]. This might create many of us nervous and feel afraid to use personalized search engines. a customized net search are not well received until it handles the privacy downside well.

2. it's very onerous to infer user info wants accurately. Users don't seem to be static. They will indiscriminately search for one thing that they're not interested in. They even look for others generally .User search histories inevitably contain noise

that is unsuitable or maybe harmful to current search. This may create personalization methods unstable.

3. Queries mustn't be handled within the same manner with respect to personalization. Personalized search may have very little result on some queries. Some work [1,]-[2]-[3] investigates whether or not current net search ranking may be comfortable for clear/unambiguous queries and so personalization is not sensible. [3] It reveals that personalized search has little result on queries with high user choice consistency. A particular personalized search conjointly has totally different effectiveness for various queries. It even hurts search accuracy below some things. For example, topical interest-based personalization, which ends up in higher performance for the query ''mouse,'' is ineffective for the question ''freemp3 transfer.'' really, relevant documents for question ''free mp3 download'' are largely classified into constant topic classes and topical interest-based personalization has no thanks to filter out desired documents. Dou et al. [3] conjointly reveal that topical interest-based personalized  search methods are tough to deploy in a very globe search engine. They improve search performance for some queries, however they'll hurt search performance for additional queries.

## 5. CONCLUSION

This paper bestowed a client-side privacy protection framework referred to as UPS for personalized net search. UPS could probably be adopted by any PWS that captures user profiles in an exceedingly ranked taxonomy. The framework allowed users to specify bespoke privacy necessities via the hierarchical profiles. additionally, UPS conjointly performed online generalization on user profiles to shield the non-public privacy while not compromising the search quality. We proposed 2 greedy algorithms, particularly GreedyDP and GreedyIL, for the net generalization. Our experimental results unconcealed that UPS may reach quality search results whereas protective user's bespoke privacy necessities. The results conjointly confirmed the effectiveness and efficiency of our resolution. For future work, we are going to attempt to resist adversaries with broader information, like richer relationship among topics (e.g., snobbery, sequentiality, and so on),or capability to capture a

series of queries (relaxing the second constraint of the person in Section three.3) from the victim. we are going to conjointly get additional subtle methodology to build the user profile, and higher metrics to predict the performance (especially the utility) of UPS.

## REFERENCES

[1] Chirita P.A., Firan C., and Nejdl W. Summarizing localcontext to personalize global web search. In Proc. Int. Conf. on Information and Knowledge Management, 2006.

[2] Chirita P.A., Nejdl W., Paiu R., and Kohlschu¨tter C. Using ODP metadata to personalize search. In Proc. 31st Annual Int. ACM SIGIR Conf. on Research and Development in Information Retrieval, 2005, pp. 178–185.

[3] Dou Z., Song R., and Wen J. A large-scale evaluation and analysis of personalized search strategies. In Proc. 33rd Annual Int. ACM SIGIR Conf. on Research and Development in Information Retrieval, 2007.

[4] Shen X., Tan B., and Zhai C. Privacy protection in personalized search. SIGIR Forum, 41(1):4–17, 2007.

[5] M. Spertta and S. Gach, "Personalizing Search Based on User Search Histories," *Proc. IEEE/WIC/ACM Int'l Conf. Web Intelligence (WI),* 2005.

[6] Y. Xu, K. Wang, B. Zhang, and Z. Chen, "Privacy-Enhancing Personalized Web Search," *Proc. 16th Int'l Conf. World Wide Web (WWW),* pp. 591-600, 2007.

[7] J.Jayanthi and Dr.K.S.Jayakumar An Integrated Page Ranking Algorithm for Personalized, Web Search, International Journal of Computer Applications (0975 – 8887), Volume 12– No.11, January 2011

[8] J. Jayanthi, M. Ezhilmathi and S. Rathi  A Novel Relevance Metric Prediction Algorithm For A Personalized Web Search,  Department of Computer Science and Engineering, Sona College of Technology, India.

[9] X. Shen, B. Tan, and C. Zhai, "Privacy Protection in Personalized Search," *SIGIR Forum,* vol. 41, no. 1, pp. 4-17, 2007.