# Providing Privacy in Profile Based Personalized Web Search

## Rajani S. Sajjan[1*], Suvarna A. Veer[2]

[1,2]Department of Computer Science & Engineering, VVPIET, Solapur University, Solapur, India

*Corresponding Author: rajanisajjan78@gmail.com*

*Abstract—* Web search engines (e.g. Google, Yahoo, Microsoft Live Search, Bing, etc.) are mostly used to search certain information from a large amount of data in a very few amount of time. Aforementioned engines are built for all kind of people and not for any particular client that is, it gives generalized result for input query and not user specific result, to address this problem personalized web search is best way to increase the accuracy of web search in terms of giving user specific results. However, effective personalized web search requires gathering and aggregating user information (e.g. user name, contact no, etc), which often raises serious concerns of privacy infringement for many users. In fact, these privacy concerns have become one of the major reasons for deploying personalized web search applications and how to do privacy-preserving personalization is a great challenge. In this proposed system, we propose and try to resist adversaries with broader background knowledge, such as richer relationship among topics. Richer relationship means we generalize the user profile results by using the background knowledge which is going to store in history. Through this we are able to hide the user search results. By using this new mechanism, we can achieve the better privacy and improve better search quality results.

*Keywords—* Data security, public server, SSM, PWS

## I. INTRODUCTION

The web search engine has most popular and important portal for common people to search certain information on the web. But sometimes, user fails due to numerous amounts of irrelevant results that do not meet user's real expectation. Such irrelevance mainly occurred due to the user's context, background and ambiguity in text [2], so that the personalized web search is best category to search useful information. But, the effective personalized web search needs gathering and aggregating user information such as user name, contact no etc. that leads to the users privacy concern[12][16].

Therefore, we propose a new protocol specially designed to provide privacy to the user's in front of web search profiling. In this proposed system, we propose and try to oppose adversaries with broader background knowledge, such as richer relationship between topics. Richer relationship means we generalize the user profile results by using the background knowledge which is going to store in history. Through this we are able to hide the user search results. In the Existing System, Greedy IL and Greedy DP algorithm are used which takes large computational time [2][10][17].

The data can be retrieved by using the background knowledge for generalization. Through this we can oppose

the adversaries. The main problem in publishing transaction data is its privacy protection. An important feature of transaction data is the extreme sparsely, which makes any single technique not sufficient in anonymizing such data. Among recent works, some suffer from high information loss, some result in data hard to interpret, and some suffer from performance drawbacks [13][19]. From some previous studies [6], it can be seen that most of the users are willing to compromise privacy if the personalization by supplying user profile to the search engine provides better search quality. In this proposed system, we propose generalization to minimize information loss. We propose new techniques to address the efficiency and scalability challenges. In the proposed System, we are going to implement the process by using which the system can become capable of capturing and extracting a series of queries by applying string similarity match algorithm to minimize the computational time and to achieve more accuracy in search results.

PWS can generally falls into two types [5], Click-log-based methods and Profile-based methods.

In Click-log-based methods we found as-

They simply impose bias to clicked pages in the user's query history.

It can only work on repeated queries from the same user, which is a strong limitation confining its applicability.

In Profile-based methods we found as-

Profile-based methods can be potentially effective for almost all sorts of queries, but are reported to be Unstable under some circumstances.

Improve the search experience with complicated user-interest models generated from user profiling techniques.

PWS has demonstrated more effective in improving the quality of web search recently, with increasing usage of personal and behavior information to profile its users, which is usually gathered implicitly from query history, browsing history, click-through data bookmarks, user documents and so forth [11][21].

## II. RELATED WORK

Authors Mrs. Sharvari V. Malthankar , Prof. Shilpa Kolte[2] implemented system a client-side privacy protection framework. System is potentially be adopted by any PWS that captures user profiles in a hierarchical taxonomy. The framework allowed users to specify customized privacy requirements via the hierarchical profiles. In addition, Online generalization on user profiles to protect the personal privacy without compromising the search quality. GreedyDP and GreedyIL algorithms are used for online generalization. Experimental results revealed that system could achieve quality search results while preserving user's customized privacy requirements.

Authors S. Manek, Aishwarya J. Reddy Vaibhavu panchal, Vijaya Pinjarkar [3] proposed an efficient information retrieval system in order to overcome the drawbacks of the ranking algorithms and improve the efficiency of web searching respecting to the precision measures. Current search engines do not rank the searched documents for a certain query automatically; they just retrieve related documents to that query issued by the user.

Authors Pratibha Rathod and Smita Desmukh [4] Proposed PMSE framework which extracts and learns user's search and location preferences based on the user's clickthrough. Here author used GPS trajectories to adapt the user mobility. Author believes that a GPS location helps to improve retrieval effectiveness, especially for location queries. Two privacy parameters, minDistance and expRatio are proposed. The privacy parameters facilitate smooth control of privacy exposure while maintaining good ranking quality

Authors Lidan Shou, He Bai, Ke Chen, and Gang Chen [5] presented a client-side privacy protection framework called UPS for personalized web search. Author believes that the proposed system could be adopted by any PWS that retrieves user's profiles in a hierarchical manner. Users can specify customized privacy requirements via the hierarchical

profiles. Proposed system also performs online generalization on user profiles to protect the personal privacy without compromising the search quality. In this paper author proposed two greedy algorithms, namely GreedyDP and GreedyIL, for the online generalization. Experimental results show that proposed system search quality results while preserving user's customized privacy requirements.

Authors Sachin S. Kale, Dattatray N. Udmale, Anjali B. Navale, Prerana S. Wagh and Prof. Rahinj P.L [6] proposed a framework for secure personalized web search. Here authors built the user profile by using domain knowledge and authors not only proposed a method to maintain the privacy and confidentiality by encrypting the user profile at the server side but also security is also provided to transportation of the data. Experimental results shows system gives better search result while using advanced user profile as compared with simple user profile on same queries.

Authors Brahmaji Katragadda, S.k. Meera [7] presented a client-side privacy protection framework. Proposed technique aims at maintaining balance between two predictive metrics that targets the utility of personalization and the privacy risk of uncovering the generalized profile, for runtime generalization author used GreedyDP and GreedyIL algorithms. At the end author imported an online prediction mechanism which decides whether personalizing a query is serviceable. Evaluated results shows that the effectiveness of proposed framework also reveals that GreedyIL surpasses GreedyDP in terms of efficiency.

Authors Anoj Kumar, Mohd. Ashraf [8] explained a promising technique that will guide the researchers to develop personalized web search technique. On studying different concepts and techniques, authors recommended collaborative filtering technique for user personalization and k-nearest neighbor partition clustering because it's have various advantages over other methods. Proposed paper mainly focuses on different techniques for efficient personalized web search and also stated merits and demerits of various available techniques.

Authors K R Remesh Babua, Philip Samuel b [9] proposes a system that uses concept network and genetic algorithm to improve the efficiency of search process. Here to store users profile author used concept network and to compare user's interest in order to predict his interest author introduced genetic algorithm. Construction of concept network is depending on the extracted pages searched by the user. Genetic Algorithm calculates the similarities between different concept network and merge these concept networks to the user's concept network and finally to the user profile. In this paper authors used TF-IDF value for extracting the concepts and creating the concept network for efficient

personalized search. In future works author will investigate the possibilities of other methods to extract the concepts from the pages browsed by the user. Also we will come with meta heuristic algorithms like Ant Colony Optimization, Particle Swarm Optimization, etc., for getting better results.

Authors Kamlesh Makvana, Pinal Shah and Parth Shah [10], introduced a novel approach to personalize   web search results by reformulating user's ambiguous  query and re-ranking algorithm. First author introduced an approach that identifies and removes the ambiguities from user's query by appending some useful keywords. System also provide related search that better identifies the current interest of user. Finally we have introduced an algorithm that re-orders the user's search result based on their preferences. Dwell time and actual rank of link is useful to incorporate user's preferences. Result analysis shown that proposed approach completely personalize search result based on user context. It has been also shown that proposed approach displayed most relevant link at top of the retrieved result.

Section I contains the introduction of personalized web search and its practices , Section II contains a survey on the related work and existing frameworks, Section III contains detail description of proposed system and architecture, Section IV contain result and analysis of proposed framework, section V contains advantages and disadvantages of proposed methodology, Section VI concludes research work with future directions).

### III.   PROPOSED SYSTEM AND ARCHITECTURE

In the proposed system, we propose a new protocol specially designed to provide privacy to the user's in front of web search profiling. In this proposed system, we propose and try to oppose adversaries with broader background knowledge, such as richer relationship among topics. We have generalized the user profile results by using the background knowledge which is going to store in history. Through this we are able to hide the user search results. In the existing system, Greedy IL and Greedy DP algorithm are used which takes large computational time [10] [15] [17].
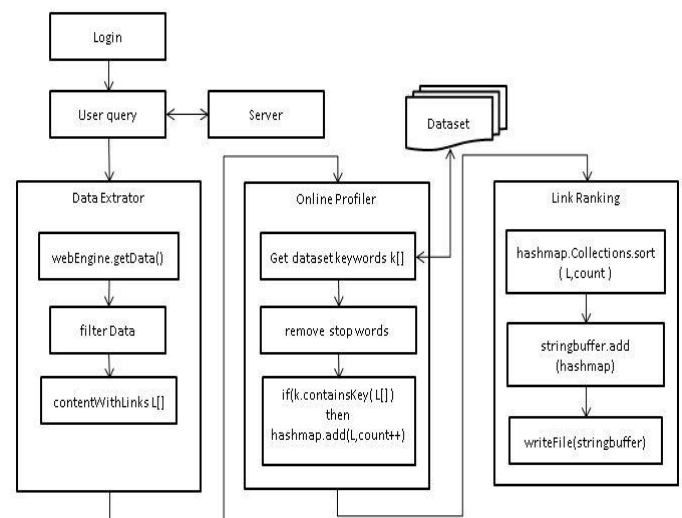


Figure 1. Proposed System Architecture

#### A. User Login
This module is for user login page. In this module, users are entered by using the user email id and password. In this module, users are entered after registering. After registering each user has unique user id and password. After login, user can posts some queries which are based on the data set which is loaded into the database.

#### B. Dataset Pre-processing
Data set is nothing but a single statistical data matrix or a single database table. Mostly content of a data set corresponds to the contents of a single statistical data matrix, or a single database table, where each column of the table shows a variable and each row is co-related to a given member of the data set in question. The data set combines lists values for each of the variables, such as height and width of an object, for each and every member of the dataset. Each value is called as a datum. The data set may combine data for one or more members, corresponding to the number of rows. Basically these modules select input dataset from registered users. Then selected dataset has been loaded into the database.  After loading the data set into the database, we are able to view the data set. Then by using the string similarity match algorithm, we filter out wanted values in the dataset and it has been pre-processed and store into the database.

#### C. Query Searching and Search Results Retrieval
In this module, user posts some queries. Depending on the submitted query, relevant results has been shown and also based on the submitted query some history results are displayed. From on the submitted query and already posted queries, we can calculate the similarity values between them. From that, the result is retrieved which is based on the more relevant results by using the maximum range of similar values.

*D. Estimate Relevant Results*

In this module, user submits query and sub query also. Based on the posted query and sub query, calculate the results based on string similarity match algorithm. Based on the relevant results and total number of data in the dataset, we can estimate the support values [14].

*E. Retrieve User Profile in Privacy Manner*

In this module, adversaries to hide the history results means, only query time has been displayed. In this module, other information such as query, query results, user name, user password are not displayed by using the background knowledge. First we generalize the table and then suppress the values based on the generalized table. Generalized values are stored in the history results. When the adversaries views the history result means, they can only view the generalized results. Finally, the performance can be evaluated by using the parameter such as computational time [18][22].

*F. Ranking*

Ranking Algorithm calculates the mean of whole links, creates a mean value and it takes nearest value from the mean value and generates the output based on score. We have used HashMap to store selected links and their respective weight. HashMap doesn't preserve any order by default. We can sort it explicitly based on the requirement.



Figure. 2 User profile tree structure

## IV. RESULT AND ANALYSIS

The performance analysis is evaluated to prove the effectiveness of the proposed methodology in terms of the comparison with the existing system [5]. Findings clearly show that the proposed system improved in retrieving the user search content according to the user's environment.

To evaluate experimental results we used AOL dataset and My dataset that we have created from search history. Below results are evaluated on AOL dataset and My dataset.

1) **My dataset:**

This dataset is created by the user itself. The query posted by user stored in the dataset with his unique id. This collection consists of 500 queries. The data is stored with user id, query and query count [1].

2) **AOL dataset:**

AOL dataset contains collection of 20M web queries collected from 650k users over three months. The data is sorted by anonymous user ID and is arranged sequentially [20].

The goal of this dataset is to provide real web query data that is based on real users. This dataset can be used for personalization, query reformulation or other types of search research.

The data set includes {AnonID, Query, QueryTime, ItemRank, ClickURL}. Where,

AnonID - an anonymous user ID number.

Query - the query issued by the user, here query is filtered with most punctuation removed.

QueryTime – Represents the time at which the query was submitted for search.

ItemRank - if the user clicked on a search result, the rank of the item on which they clicked is listed.

ClickURL - if the user clicked on a search result, the domain portion of the URL in the clicked result is listed [20].

We have filtered AOL dataset and minimized the parameters into UID, Query and Query count.
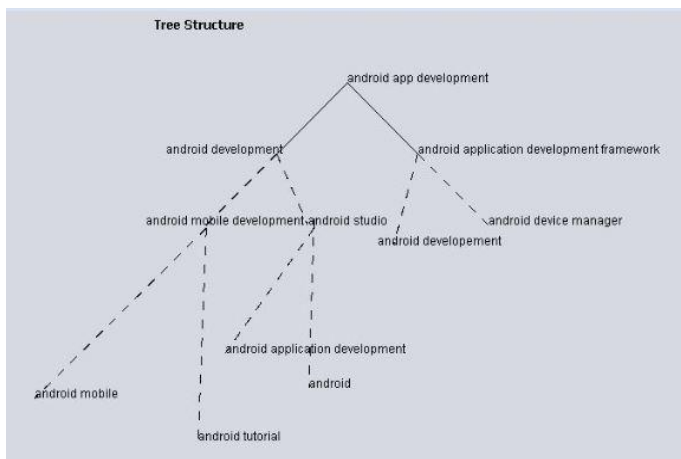
Table1.System modules

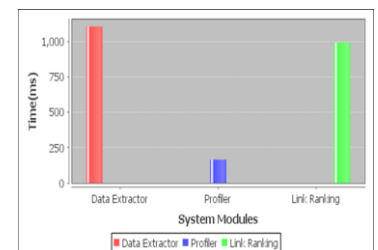| System Modules | Time (ms) |
|---|---|
| Data Extractor | 1120 |
| Profiler | 198 |
| Link Ranking | 991 |



Figure 3. System Modules

Fig 3 shows the time required for specific modules, here data extractor time is depend on the internet speed, and for profiler and link ranking module increase in data is directly proportional to increase in time.

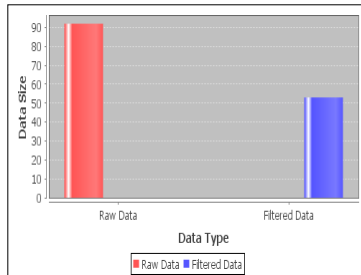| Data Type | Data Size |
|-----------|-----------|
| Raw Data | 92 |
| Filtered data | 52 |

Table 2. Data Size



Figure. 4 Data Size

Fig 4 shows the raw data extracted by extractor from search engine and filtered data filtered by profiler depending on the user's profile, proposed system minimize data by on an average by 30 percent.

Table 3. System timings

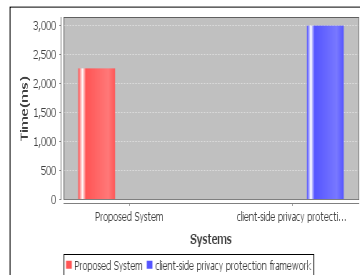| Systems | Time (ms) |
|---------|-----------|
| Proposed System | 2290 |
| client-side privacy protection framework | 3000 |



Figure. 5 System timings

Fig 5 shows the total execution time required for proposed system and existing system that is client side privacy protection framework. Graph shows proposed system take minimum time to execute as compared to existing.

Table 4. System timings on variable dataset size

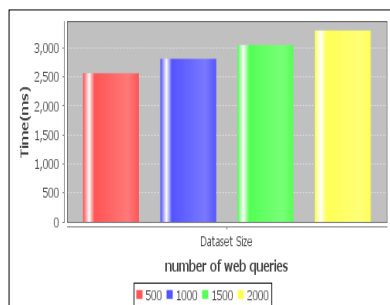| No. of web queries | Time (ms) |
|--------------------|-----------|
| 500 | 2500 |
| 1000 | 2710 |
| 1500 | 3100 |
| 2000 | 3200 |



Fig. 6 System timings on variable dataset size

Fig 6 shows the total execution time required for proposed system on variable size dataset. Graph clearly depicts that increase in data size is directly proportional to increase in execution timing.

**In proposed system performance is evaluated with 500 web queries.**

1) **Accuracy:** It is the degree to which the result of a measurement, calculation, or specification conforms to the correct value (true). The proposed system gives the 96.6% Accuracy.

$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN} *100$$

Where,
TP (True Positive) – Correctly Identified. = 21
FP (False Positive) – Incorrectly Identified. = 9
TN (True Negative) – Correctly Rejected.= 462
FN (False Negative) – Incorrectly Rejected. = 8

2) **Sensitivity:** It is a True Positive Rate (TPR). It measures propagation of actual positives which are correctly identified. It measures the percentage of correctly identified. The proposed system gives 72.41% Sensitivity

$$Sensitivity = \frac{TP}{TP+FN} *100$$

3) **Specificity:** It is True Negative Rate. It measures propagation of negatives which are correctly identified. It measures the percentage of correctly identified. The proposed system gives 98.08 % Specificity

$$Specificity = \frac{TN}{TN+FP} *100$$

4) **Precision:** It is degree to which repeated measurements under unchanged condition show the same results. The proposed system gives 70 precision

$$Precision = \frac{TP}{TP+FP} *100$$

Table 5.System performance

| Parameters | Result % |
|------------|----------|
| Accuracy | 96.6 |
| Sensitivity | 72.41 |
| Specifity | 98.08 |
| Precision | 70 |



Fig. 7 System performance

Table 6. System Accuracy

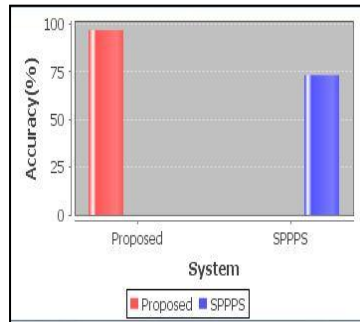| Framework | Accuracy |
|---|---|
| Proposed | 96.6 |
| Supporting Privacy Protection in Personalized Web Search | 73.2 |



Fig. 8 System Accuracy

Table 7. Comparative Analysis

| Parameters | Google [2] | Existing framework [5] | Proposed framework |
|---|---|---|---|
| Overall Precision | 27 % | 66.2% | 70% |
| Integration | It's a standalone search engine. | It's a standalone framework cannot integrate with other frameworks | can be used on top of any search engine |
| Searching algorithm | Google use PageRank algorithm. It counts the number and quality of links in a page to determine a rough estimate of how important the website is. | Use greedyDP and greedyIL algorithm which works on words rather than page | Use string similarity algorithms and word weighting to determine important word from profile that user has searched before and rerank links according to word weight |
| General working | Google is widely used to find certain data among a huge amount of information in a minimal amount of time. However, these search tools pose a privacy threat to the users: web search engines profile their users by storing and analyzing past searches submitted by them | Runtime generalization focuses on maintaining balance between two predictive metrics that evaluate the utility of personalization and the privacy risk of exposing the generalized profile. | Proposed framework can adaptively generalize profiles by queries while respecting user specified privacy requirements The key component for privacy protection is an online profiler implemented as a search proxy running on the client machine itself. |

Above table 7 shows the comparative analysis between universal search engine Google, existing framework "Supporting Privacy Protection in Personalized Web Search" by Lidan Shou, He Bai, Ke Chen, and Gang Chen and our proposed framework

## V. ADVANTAGES AND DISADVANTAGES

1) Proposed framework can adaptively generalize profiles by queries while respecting user specified privacy requirements.
2) It provides better search results with accuracy around 96.9 percent while the accuracy of existing system is 73.2[5].
3) It has less computational time of 2290 millisecond as compared to existing system 3000 millisecond[5].
4) More scalable in terms of computation complexity

## VI. CONCLUSION AND FUTURE SCOPE

The data can be retrieved by using the background knowledge for generalization. An important feature of transaction data is the extreme sparsity, which makes any single technique not sufficient in anonymizing such data. Among recent works, some suffer from high information loss, some result in data hard to interpret, and some suffer from performance drawbacks. From some previous studies, it can be seen that most of the users are willing to compromise privacy if the personalization by supplying user profile to the

search engine provides better search quality. In the proposed system, we propose generalization to minimize information loss. We propose new techniques to address the efficiency and scalability challenges with better accuracy. An evaluated result shows that system surpassed the accuracy of existing system with better precision rate of 70%. In the future we would try to enhance the search quality and provide more security from the adversaries.

## REFERENCES

[1] Ms. Suvarna A Veer, Rajani S. Sajjan Computer Science & Engineering Department, VVPIET, Solapur, India *"Providing Privacy in Domain Specific Search with SSM and Cosine Similarity"* www.jetir.org Volume 5, Issue 12 JETIR December2018,

[2] Mrs. Sharvari V. Malthankar , Prof. Shilpa Kolte PG Student, *"Client side Privacy Protection Using Personalized Web Search"* Elsevier Science direct 7th International Conference on Communication, Computing and Virtualization 2016

[3] S. Manek fmanek3@gmail.com Aishwarya J. Reddy Vaibhavu panchal Vijaya Pinjarkar Department of Information Technology K.J.Somaiya Institute of Engineering and Information Technology, Sion. Mumbai, Maharashtra, India vkhirodkar@somaiya.edu *"Hybrid Crawling for Time-Based Personalized Web Search Ranking Foram"* 978-1-5090-5686-6/17/$31.00 ©2017 IEEE

[4] Pratibha Rathod pratima.rathod15@gmail.com Smita Desmukh Information Technology, deshmukhsmita17@yahoo.com *"A Personalized Mobile Search Engine based on User Preference"* IEEE International Conference on Power, Control, Signals and Instrumentation Engineering (ICPCSI-2017)

[5] Lidan Shou, He Bai, Ke Chen, and Gang Chen,*"Supporting PrivacyProtection in Personalized Web Search",*IEEE Transactions on Knowledge and Data engineering,vol.26,no.2,february 2014

[6] Sachin S. Kale1, Dattatray N. Udmale1, Anjali B. Navale1, Prerana S. Wagh1, Prof. Rahinj P.L2 B.E *"Supporting Privacy Protection in Personalized Web Search, India International Journal of Innovative Research in Computer and Communication Engineering* 2017

[7] Brahmaji Katragadda, 2Sk.Meera *"Supporting Privacy Protection in Personalized Web Search"* Journal of Science and Technology (JST) Volume 2, Issue 7, July 2017

[8] Anoj Kumar anoj.kr@hotmail.com Mohd. Ashraf ashraf.saifee@gmail.com *"Personalized Web Search Engine using Dynamic User Profile and Clustering Techniques"* 978- 9-3 805-4416-8/15/$31. 00 c2 01 5 IEEE

[9] K R Remesh Babua,Philip Samuelb, *"Concept Networks for Personalized Web Search Using Genetic Algorithm"* International Conference on Information and Communication Technologies 2016

[10] Kamlesh Makvana, Pinal Shah, Parth Shah, kamleshmakvana.it@charusat.ac.in, parthshah.ce@charusat.ac.in pinalshah.it@charusat.ac.in *"A Novel Approach to Personalize Web Search through User Profiling and Query Reformulation"* 978-1-4799-4674-7/14/$31.00©2014 IEEE

[11] V.Ramya, S.Gowthami *"ENHANCE PRIVACY SEARCH IN WEB SEARCH ENGINE USING GREEDY ALGORITHM"* International Journal of Scientific Research Engineering & Technology (IJSRET), ISSN 2278 – 0882 Volume 3, Issue 8, November 2014

[12] Zhan Su, Byung-Ryul Ahn, Ki-yol Eom, Min-koo Kang, Jin-Pyung Kim, Moon-Kyun Kim Department of Artificial Intelligence, University of Sungkyunkwan Cheoncheon dong, Jangan-gu, Suwon, Korea *"Plagiarism Detection Using the Levenshtein Distance and Smith-Waterman Algorithm"* The 3rd Intetnational Conference on Innovative Computing Information and Control (ICICIC'08) 978-0-7695-3161-8/08 $25.00 © 2008 IEEE

[13] Alfirna Rizqi Lahitani1, Adhistya Erna Permanasari, Noor Akhmad Setiawan Department of Electrical Engineering and Information Technology, Faculty of Engineering Universitas Gadjah Mada *"Cosine Similarity to Determine Similarity Measure: Study Case in Online Essay Assessment"* ACM 2011

[14] ZHENGHUA XU1 , OANA TIFREA-MARCIUSKA2 , THOMAS LUKASIEWICZ1 ,MARIA VANINA MARTINEZ3 , GERARDO I. SIMARI3 , and CHENG CHEN *"Lightweight Tag-Aware Personalized Recommendation on the Social Web Using Ontological Similarity"* 4 2169-3536 (c) 2018 IEEE

[15] Yayuan Tang 1;2 , Hao Wang 3 , Kehua Guo 2;4 , Yizhe Xiao 2 , Tao Chi *"Relevant Feedback Based Accurate and IntelligentRetrieval on Capturing User Intention for Personalized Websites"* 2169-3536 (c) 2018 IEEE

[16] Mohammad Mustaneer Rahman, and Nor Aniza Abdullah, *"A Personalised Group-Based Recommendation Approach for Web Search in E-Learning"* IEEE 2169-3536 (c) 2018 IEEE.

[17] Puxuan Yu Wuhan University Wuhan, China pxyuwhu@gmail.com Wasi Uddin Ahmad wasiahmad@ucla.edu *"Hide-n-Seek: An Intent-aware Privacy Protection Plugin for Personalized Web Search"* 18, July 8-12, 2018, Ann Arbor, MI, USA Italy. 2016 ACM. ISBN 978-1-4503-4069-4/16/07

[18] Gerard Deepak, B. N. Shwetha, C. N. Pushpa, J. Thriveni & K. R. Venugopal *"A hybridized semantic trust-based framework for personalized web page recommendation"* International Journal of Computers and Applications ISSN: 1206-212X (Print) 1925-7074

[19] Najneen Tamboli, Sathish Kumar *"Review on Privacy Preservation in Personalized Web Search"* International Journal of Innovative Research in Computer and Communication Engineering (An ISO 3297: 2007 Certified Organization) Vol. 3, Issue 11, November 2015

[20] Avi Arampatzis, George Drosatos and Pavlos S. Efraimidis, *"A Versatile Tool for Privacy -Enhanced Web Searc"* Xant hi 67 100, Greece Springer- Verlag Berlin Heidelberg 2013

[21] B. SekharBabu, P. Lakshmi Prasanna, D. Rajeswara Rao, J. LakshmiAnusha, A. Pratyusha and A. Ravi Chand, *"PROFILE BASED PERSONALIZED WEB SEARCH USING GREEDY ALGORITHMS"* MAY 2016 ISSN 1819-6608 ARPN Journal of Engineering and Applied Sciences

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

[23] Hina Ansari Mahakal Institute of Technology, Ujjain *"Issues and Challenges in Measuring Security Threats During Personalized Web Search"* , Volume-3, Issue-6 ISSN: 2320-7639, IJSRCSE, 2015

**Authors Profile**

Ms. Sajjan R.S. received her M.Tech in Computer Science and Engineering. She has a working experience of 17 years and is currently the H.O.D. of the Computer Science and Engineering Department. She is currently pursuing Ph.D. Her research interest is in Cloud Computing.

.
Ms. Suvarna Arun Veer received Bachelor of Engineering in Computer Science from T COE Osmanabad. She had a working experi Of 3 Years. She is currently working toward M.E. degree in Computer Science & Engine From Solapur University, Solapur. Her resea Interests lies in area of programming & Data Mining.