

Mining Client Conduct to Gather High Likely Pursuit Objectives: Another Objective Mining

M Sophia

M.Tech (CSE)

Department of Computer Science and Engineering,
GIET Engineering College, Rajamahendravaram,
Andhra Pradesh 533296, India.

Y Durga Prasad, M.Tech

Assistant Professor,

Department of Computer Science and Engineering,
GIET Engineering College, Rajamahendravaram,
Andhra Pradesh 533296, India.

Abstract

For an expansive theme and questionable inquiry, diverse clients may have distinctive pursuit objectives when they submit it to a web crawler. The surmising and investigation of client seek objectives can be extremely valuable in enhancing web crawler pertinence and client encounter. In this paper, we propose a novel way to deal with derive client look objectives by dissecting internet searcher question logs. To begin with, we propose a system to find distinctive client scan objectives for an inquiry by bunching the proposed criticism sessions. Criticism sessions are built from client navigate logs and can effectively mirror the data needs of clients. Second, we propose a novel way to deal with produce pseudo-archives to better speak to the criticism sessions for grouping. At long last, we propose another model "Arranged Average Precision (CAP)" to assess the execution of deducing client look objectives. Trial comes about are exhibited utilizing client navigate logs from a business web index to approve the adequacy of our proposed techniques.

Introduction to Data Mining:

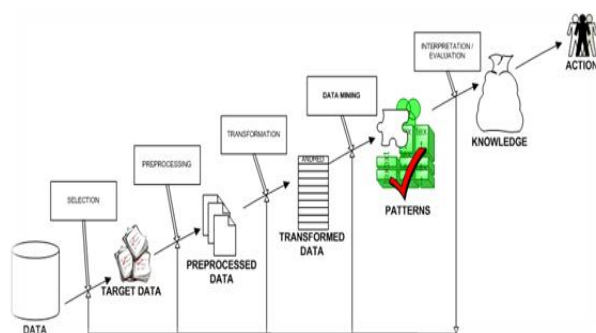


Fig 1.1 Structure of Data Mining

Data mining (sometimes called data or knowledge discovery) is the process of analyzing data from different perspectives and summarizing it into useful information - information that can be used to increase revenue, cuts costs, or both. Data mining software is one of a number of analytical tools for analyzing data. It allows users to analyze data from many different dimensions or angles, categorize it, and summarize the relationships identified. Technically, data mining is the process of finding correlations or patterns among dozens of fields in large relational databases [1-4].

Working of Data Mining

While large-scale information technology has been evolving separate transaction and analytical systems, data mining provides the link between the two. Data mining software analyzes relationships and patterns in stored transaction data based on open-ended user queries. Several types of analytical software [6] are available: statistical, machine learning, and neural networks.

Data mining consists of five major elements:

1. Extract, transform, and load transaction data onto the data warehouse system.
2. Store and manage the data in a multidimensional database system.
3. Provide data access to business analysts and information technology professionals.
4. Analyze the data by application software.

Cite this article as: M Sophia & Y Durga Prasad, "Mining Client Conduct to Gather High Likely Pursuit Objectives: Another Objective Mining", International Journal of Research in Advanced Computer Science Engineering, Volume 4 Issue 1, 2018, Page 1-9.

5. Present the data in a useful format, such as a graph or table.

Characteristics of Data Mining:

Large quantities of data:

The volume of data is so great it has to be analyzed by automated techniques e.g. satellite information, credit card transactions etc.

Noisy, incomplete data:

Imprecise data is the characteristic of all data collection.
Complex data structure: conventional statistical analysis not possible

Heterogeneous data stored in legacy systems

Benefits of Data Mining:

- 1) It's one of the most effective services that are available today. With the help of data mining, one can discover precious information about the customers and their behavior for a specific set of products and evaluate and analyze, store, mine and load data related to them
- 2) An analytical CRM model [8] and strategic business related decisions can be made with the help of data mining as it helps in providing a complete synopsis of customers
- 3) An endless number of organizations have installed data mining projects and it has helped them see their own companies make an unprecedented improvement in their marketing strategies (Campaigns)
- 4) Data mining is generally used by organizations with a solid customer focus. For its flexible nature as far as applicability is concerned is being used vehemently in applications to foresee crucial data including industry analysis and consumer buying behaviors
- 5) Fast paced and prompt access to data along with economic processing techniques have made data mining one of the most suitable services that a company seeks

EXISTING SYSTEM

We define user search goals as the information on different aspects of a query that user groups want to obtain. Information need is a user's particular desire to obtain information to satisfy his/her need. User search goals can be considered as the clusters of information

needs for a query. The inference and analysis of user search goals can have a lot of advantages in improving search engine relevance and user experience [10].

DISADVANTAGES OF EXISTING SYSTEM

- What users care about varies a lot for different queries, finding suitable predefined search goal classes is very difficult and impractical.
- Analyzing the clicked URLs directly from user click-through logs to organize search results. However, this method has limitations since the number of different clicked URLs of a query may be small. Since user feedback is not considered, many noisy search results that are not clicked by any users may be analyzed as well. Therefore, this kind of methods cannot infer user search goals precisely.
- Only identifies whether a pair of queries belongs to the same goal or mission and does not care what the goal is in detail.

PROPOSED SYSTEM

In this paper, we aim at discovering the number of diverse user search goals for a query and depicting each goal with some keywords automatically [12]. We first propose a novel approach to infer user search goals for a query by clustering our proposed feedback sessions. Then, we propose a novel optimization method to map feedback sessions to pseudo-documents which can efficiently reflect user information needs. At last, we cluster these pseudo documents to infer user search goals and depict them with some keywords.

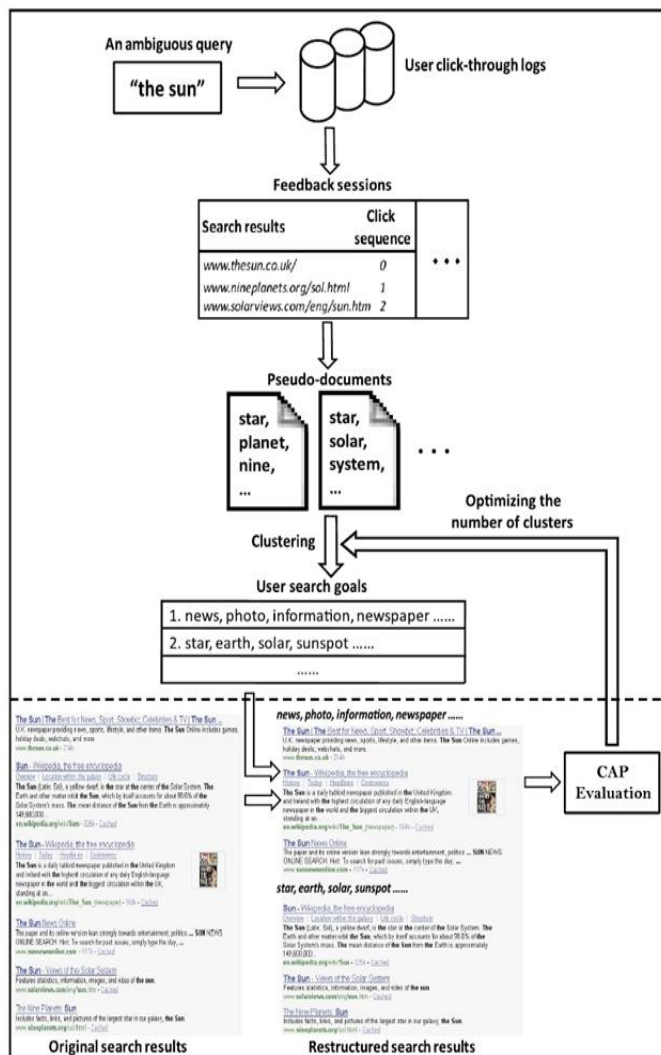
ADVANTAGES OF PROPOSED SYSTEM

- To sum up, our work has three major contributions as follows:
- We propose a framework to infer different user search goals for a query by clustering feedback sessions. We demonstrate that clustering feedback sessions is more efficient than clustering search results or clicked URLs directly. Moreover, the distributions of different

user search goals can be obtained conveniently after feedback sessions are clustered.

- We propose a novel optimization method to combine the enriched URLs in a feedback session to form a pseudo-document, which can effectively reflect the information need of a user. Thus, we can tell what the user search goals are in detail.
- We propose a new criterion CAP to evaluate the performance of user search goal inference based on restructuring web search results. Thus, we can determine the number of user search goals for a query.

SYSTEM ARCHITECTURE:



IMPLEMENTATION MODULES

- Feedback Sessions
- Pseudo-documents
- Inferring Pseudo-documents
- Evaluation Search Result

MODULES DESCRIPTION

Feedback Sessions

The inferring user search goals for a particular query. Therefore, the single session containing only one query is introduced, which distinguishes from the conventional session. Meanwhile, the feedback session in this paper is based on a single session, although it can be extended to the whole session [14]. The proposed feedback session consists of both clicked and unclicked URLs and ends with the last URL that was clicked in a single session. It is motivated that before the last click, all the URLs have been scanned and evaluated by users. Therefore, besides the clicked URLs, the unclicked ones before the last click should be a part of the user feedbacks.

Pseudo-documents

The URLs with additional textual contents by extracting the titles and snippets of the returned URLs appearing in the feedback session. In this way, each URL in a feedback session is represented by a small text paragraph that consists of its title and snippet. Then, some textual processes are implemented to those text paragraphs, such as transforming all the letters to lowercases, stemming and removing stop words. to obtain the feature representation of a feedback session, we propose an optimization method to combine both clicked and unclicked URLs in the feedback session.

Inferring Pseudo-documents

The proposed pseudo-documents, we can infer user search goals. In this section, we will describe how to infer user search goals and depict them with some meaningful keywords. As each feedback session is represented by a pseudo-document and the feature representation of the pseudo-document. pseudo-documents by K-means clustering [16] which is simple

and effective. Since we do not know the exact number of user search goals for each query, we set K to be five different values and perform clustering based on these five values, respectively. The terms with the highest values in the center points are used as the keywords to depict user search goals. Note that an additional advantage of using this keyword based description is that the extracted keywords can also be utilized to form a more meaningful query in query recommendation and thus can represent user information needs more effectively.

Evaluation Search Result

If user search goals are inferred properly, the search results can also be restructured properly, since restructuring web search results is one application of inferring user search goals. Therefore, we propose an evaluation method based on restructuring web search results to evaluate whether user search goals are inferred properly or not. In this section, we propose this novel criterion "Classified Average Precision" to evaluate the restructure results. Based on the proposed criterion, we also describe the method to select the best cluster number [5].

ALGORITHMS

CAP

Capturing Feedback Sessions:

For a web search is a series of successive queries to satisfy a single information need and some clicked search results. Here, feedback session consists of both clicked and unclicked URL's and ends with the last URL that was clicked in a single session. Clicked URL's state what users require and unclicked URL's reflect what users do not care about. For inferring user search goals it is more efficient to analyze the feedback sessions than to analyze search results or clicked URL's directly because there are different feedback sessions in user click-through logs [7].

Building pseudo-documents:

Representing the URL's in feedback session. Each URL's title and snippet are represented by term frequency-inverse document frequency as below,

$$\begin{aligned} \mathbf{T}_{u_i} &= [t_{w_1}, t_{w_2}, \dots, t_{w_n}]^T, \\ \mathbf{S}_{u_i} &= [s_{w_1}, s_{w_2}, \dots, s_{w_n}]^T, \end{aligned} \quad (1)$$

Where \mathbf{T}_{u_i} and \mathbf{S}_{u_i} are TF-IDF vectors of the URL's title snippet. u_i means i th URL in the feedback session.

$w_j(j=1,2,\dots,n)$ is j th term appearing in the enriched URL.

$$\mathbf{F}_{u_i} = \omega_t \mathbf{T}_{u_i} + \omega_s \mathbf{S}_{u_i} = [f_{w_1}, f_{w_2}, \dots, f_{w_n}]^T, \quad (2)$$

Here, \mathbf{F}_{u_i} is feature representation of i th URL in feedback session. ω_t and ω_s are weights of title and snippet. Here title should be more significant than snippets. So, the weight of title should be higher.

Forming pseudo-documents based on URL representations:

Here, an optimization method is used to combine both clicked and unclicked URL's in the feedback sessions. Let \mathbf{F}_{fs} be the feature representation of feedback sessions and $f_{fs}(w)$ be the value for term w . $\mathbf{F}_{ucm}(m=1,2,\dots,M)$ and $\mathbf{F}_{ucl}(l=1,2,\dots,L)$ be the representation of clicked and unclicked URL's in the feedback sessions. $f_{ucm}(w)$ and $f_{ucl}(w)$ are the values of term w in vectors. Obtain such a \mathbf{F}_{fs} that sum of distances between \mathbf{F}_{fs} and each \mathbf{F}_{ucm} is minimized and sum of distances between \mathbf{F}_{fs} and each \mathbf{F}_{ucl} is maximized [9].

$$\begin{aligned} \mathbf{F}_{fs} &= [f_{fs}(w_1), f_{fs}(w_2), \dots, f_{fs}(w_n)]^T, \\ f_{fs}(w) &= \arg \min_{f_{fs}(w)} \left\{ \sum_M [f_{fs}(w) - f_{ucm}(w)]^2 \right. \\ &\quad \left. - \lambda \sum_L [f_{fs}(w) - f_{ucl}(w)]^2 \right\}, f_{fs}(w) \in I_c. \end{aligned} \quad (3)$$

Let I_c be the interval $[\mu f_{uc}(w) - \sigma f_{uc}(w), \mu f_{uc}(w) + \sigma f_{uc}(w)]$ and I_{c^-} be the interval $[\mu f_{uc^-}(w) - \sigma f_{uc^-}(w), \mu f_{uc^-}(w) + \sigma f_{uc^-}(w)]$, where $\mu f_{uc}(w)$ and $\sigma f_{uc}(w)$ represent the mean and mean square error of $f_{uc}(w)$ respectively, and $\mu f_{uc^-}(w)$ and $\sigma f_{uc^-}(w)$ represent the mean and mean square error of $f_{uc^-}(w)$, respectively. If $I_c \in I_{c^-}$ or $I_{c^-} \in I_c$, we consider that the user does not care about the term w . In this situation, we set $f_{fs}(w)$ to be 0, as shown in

$$f_{fs}(w) = 0, I_c \subseteq I_{c^-} \text{ or } I_{c^-} \subseteq I_c. \quad (4)$$

As in (3) and (4), each feedback session is represented by a pseudo-document and the feature representation of the

pseudo-document is F_{fs} . The similarity between two pseudo-documents is computed as the cosinescore of F_{fs_i} and F_{fs_j} , as follows:

$$\begin{aligned} Sim_{i,j} &= \cos(F_{fs_i}, F_{fs_j}) \\ &= \frac{F_{fs_i} \cdot F_{fs_j}}{|F_{fs_i}| |F_{fs_j}|} \end{aligned} \quad (5)$$

And the distance between two feedback session is

$$Dis_{i,j} = 1 - Sim_{i,j}. \quad (6)$$

We cluster pseudo-documents by K-means clustering which is simple and effective. Since we do not know the exact number of user search goals for each query, we set K to be five different values (i.e., 1; 2; . . . ; 5) and perform clustering based on these five values, respectively. The optimal value will be determined through the evaluation criterion presented.

After clustering all the pseudo-documents, each cluster can be considered as one user search goal. The center point of a cluster is computed as the average of the vectors of all the pseudo-documents in the cluster, as shown in

$$F_{center_i} = \frac{\sum_{k=1}^{C_i} F_{fs_k}}{C_i}, (F_{fs_k} \in Cluster\ i), \quad (7)$$

where F_{center_i} is the i th cluster's center and C_i is the number of the pseudo-documents in the i th cluster. F_{center_i} is utilized to conclude the search goal of the i th cluster.

In order to apply the evaluation method to large-scale data, the single sessions in user click-through logs are used to minimize manual work. Because from user click-through logs, we can get implicit relevance feedbacks, namely "clicked" means relevant and "unclicked" means irrelevant [11-13]. A possible evaluation criterion is the average precision (AP) which evaluates according to user implicit feedbacks. AP is the average of precisions computed at the point of each relevant document in the ranked sequence, as shown in

$$AP = \frac{1}{N^+} \sum_{r=1}^N rel(r) \frac{R_r}{r}, \quad (8)$$

where N^+ is the number of relevant (or clicked) documents in the retrieved ones, r is the rank, N is the

total number of retrieved documents, $rel()$ is a binary function on the relevance of a given rank, and R_r is the number of relevant retrieved documents of rank r or less.

VAP is still an unsatisfactory criterion. Considering an extreme case, if each URL in the click session is categorized into one class, VAP will always be the highest value namely 1 no matter whether users have so many search goals or not. Therefore, there should be a risk to avoid classifying search results into too many classes by error. We propose the risk as follows:

$$Risk = \frac{\sum_{i,j=1(i < j)}^m d_{ij}}{C_m^2}. \quad (9)$$

It calculates the normalized number of clicked URL pairs that are not in the same class, where m is the number of the clicked URLs. If the pair of the i th clicked URL and the j th clicked URL are not categorized into one class, d_{ij} will be 1; otherwise, it will be 0. $C_m^2 = (m(m-1)/2)$ is the total number of the clicked URL pairs.

Based on the above discussions, we can further extend VAP by introducing the above Risk and propose a new criterion "Classified AP," as shown below

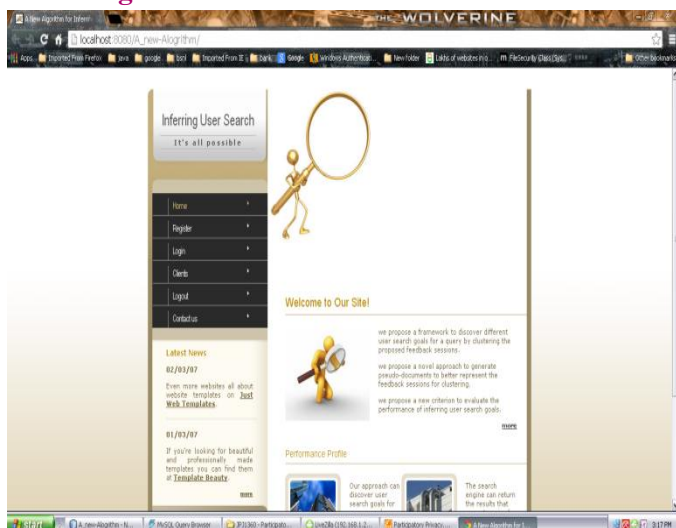
$$CAP = VAP \times (1 - Risk)^\gamma. \quad (10)$$

From (10), we can see that CAP selects the AP of the class that user is interested in (i.e., with the most clicks/votes) and takes the risk of wrong classification into account. And γ is used to adjust the influence of Risk on CAP, which can be learned from training data. Finally, we utilize CAP to evaluate the performance of restructuring search results.

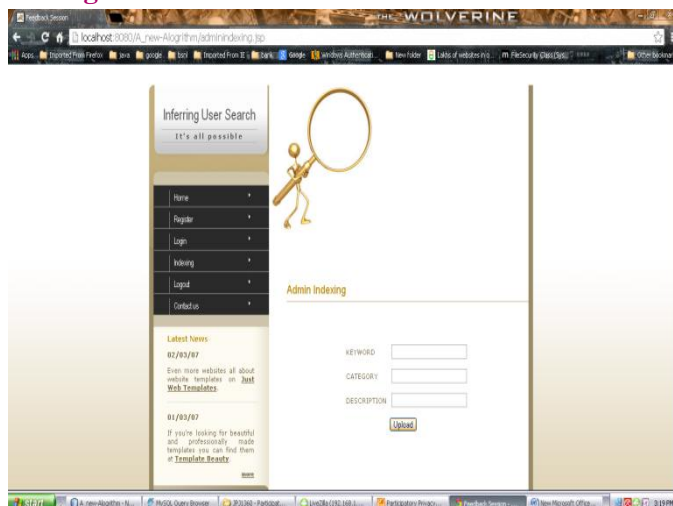
Considering another extreme case, if all the URLs in the search results are categorized into one class, Risk will always be the lowest namely 0; however, VAP could be very low. Generally, categorizing search results into less clusters will induce smaller Risk and bigger VAP, and more clusters will result in bigger Risk and smaller VAP. The proposed CAP depends on both of Risk and VAP [15].

SCREEN SHOTS

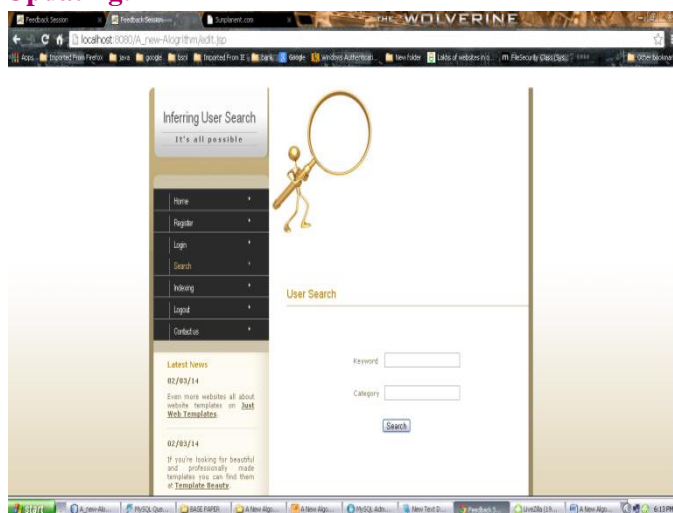
Home Page:



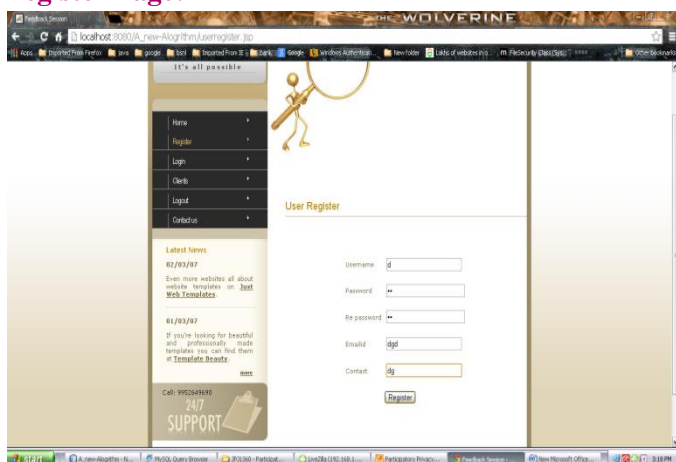
Adding information:



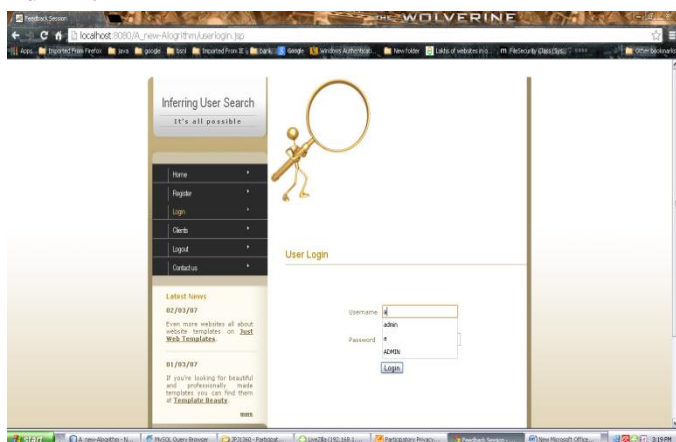
Updating:



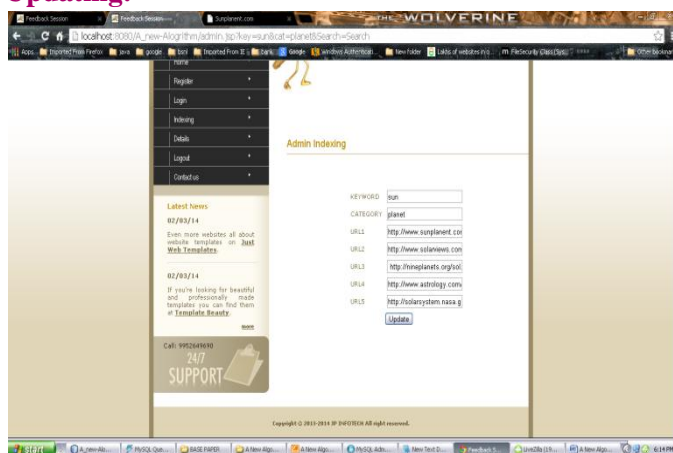
Register Page:



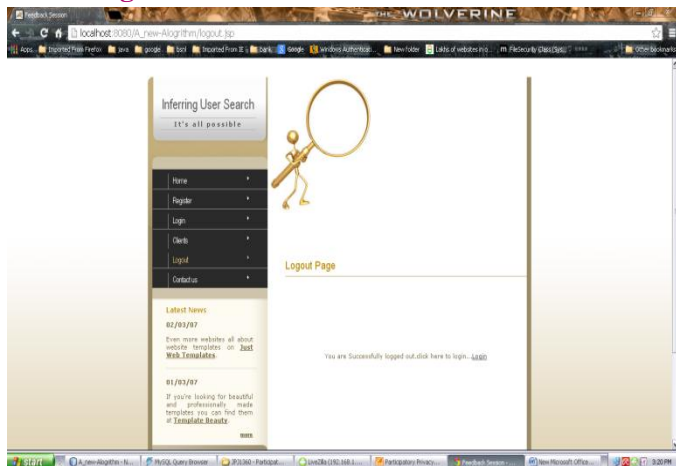
Admin:



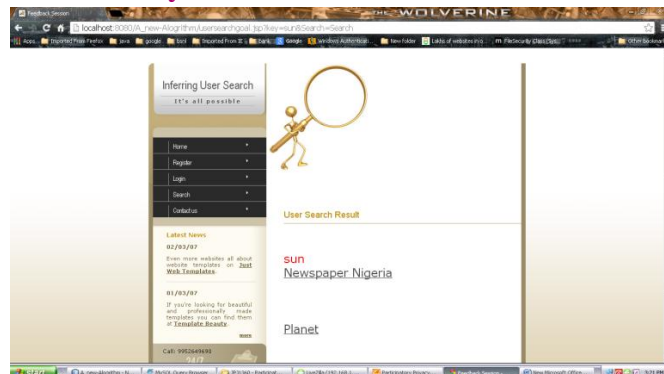
Updating:



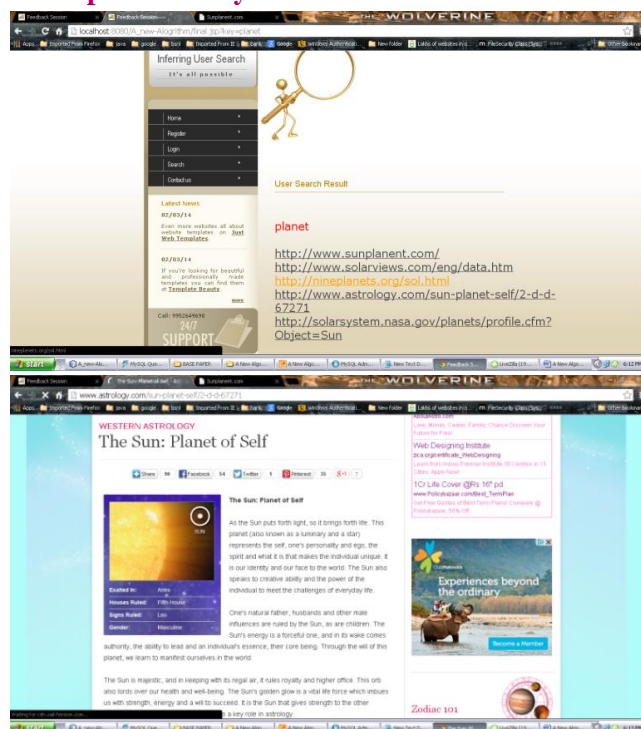
Admin logout:



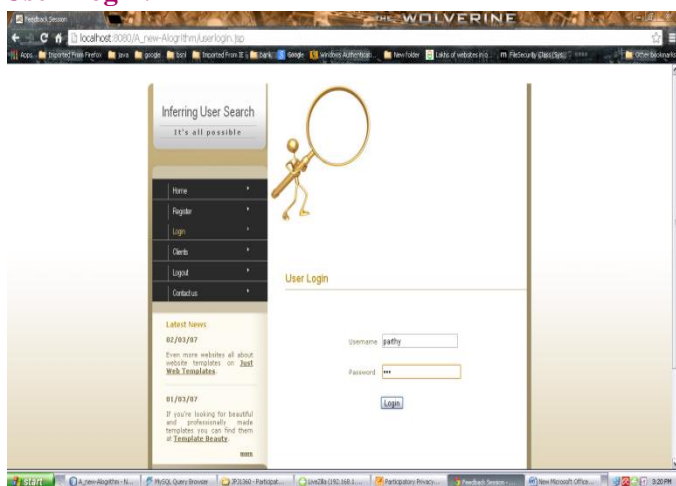
Result for keyword:



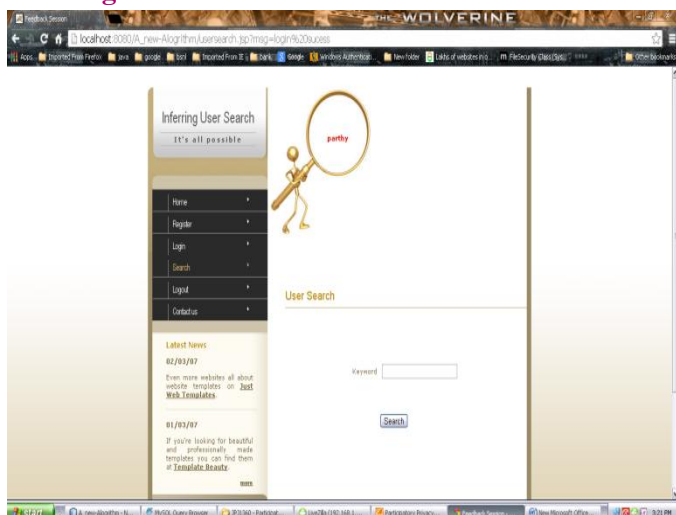
Description for keyword:



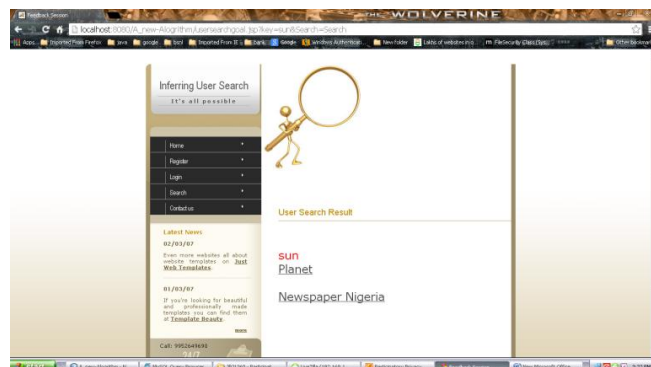
User Login:



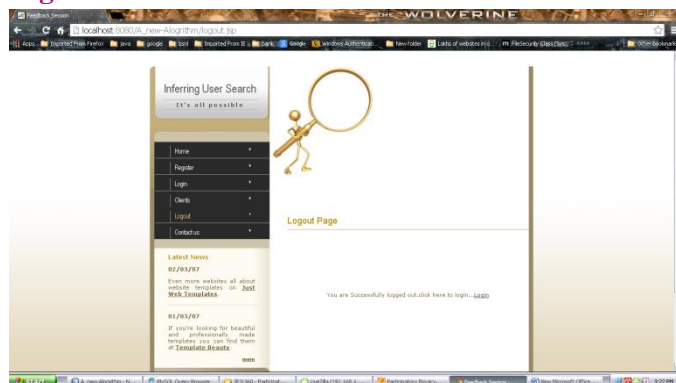
User Page:



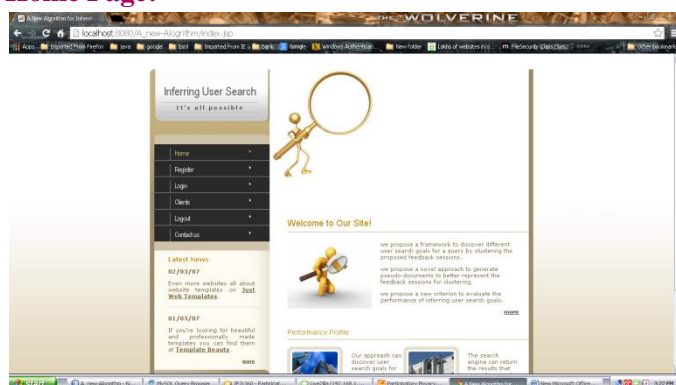
Final result:



Logout:



Home Page:



CONCLUSION

In this project, a novel approach has been proposed to infer user search goals for a query by clustering its feedback sessions represented by pseudo-documents. First, we introduce feedback sessions to be analyzed to infer user search goals rather than search results or clicked URLs. Both the clicked URLs and the unclicked ones before the last click are considered as user implicit feedbacks and taken into account to construct feedback sessions. Therefore, feedback sessions can reflect user information needs more efficiently. Second, we map feedback sessions to pseudo documents to approximate goal texts in user minds. The pseudo-documents can enrich the URLs with additional textual contents including the titles and snippets. Based on these pseudo-documents, user search goals can then be discovered and depicted with some keywords. Finally, a new criterion CAP is formulated to evaluate the performance of user search goal inference. Experimental results on user

click-through logs from a commercial search engine demonstrate the effectiveness of our proposed methods. The complexity of our approach is low and our approach can be used in reality easily. For each query, the running time depends on the number of feedback sessions. However, the dimension of Ffs in (3) and (5) is not very high. Therefore, the running time is usually short. In reality, our approach can discover user search goals for some popular queries offline at first. Then, when users submit one of the queries, the search engine can return the results that are categorized into different groups according to user search goals online. Thus, users can find what they want conveniently.

References

- [1] R. Baeza-Yates and B. Ribeiro-Neto, Modern Information Retrieval. ACM Press, 1999.
- [2] R. Baeza-Yates, C. Hurtado, and M. Mendoza, "Query Recommendation Using Query Logs in Search Engines," Proc. Int'l Conf. Current Trends in Database Technology (EDBT '04), pp. 588-596, 2004.
- [3] D. Beeferman and A. Berger, "Agglomerative Clustering of a Search Engine Query Log," Proc. Sixth ACM SIGKDD Int'l Conf. Knowledge Discovery and Data Mining (SIGKDD '00), pp. 407-416, 2000.
- [4] S. Beitzel, E. Jensen, A. Chowdhury, and O. Frieder, "Varying Approaches to Topical Web Query Classification," Proc. 30th Ann. Int'l ACM SIGIR Conf. Research and Development (SIGIR '07), pp. 783-784, 2007.
- [5] H. Cao, D. Jiang, J. Pei, Q. He, Z. Liao, E. Chen, and H. Li, "Context-Aware Query Suggestion by Mining Click-Through," Proc. 14th ACM SIGKDD Int'l Conf. Knowledge Discovery and Data Mining (SIGKDD '08), pp. 875-883, 2008.
- [6] H. Chen and S. Dumais, "Bringing Order to the Web: Automatically Categorizing Search Results," Proc.

SIGCHI Conf. Human Factors in Computing Systems (SIGCHI '00), pp. 145-152, 2000.

[7] C.-K Huang, L.-F Chien, and Y.-J Oyang, "Relevant Term Suggestion in Interactive Web Search Based on Contextual Information in Query Session Logs," J. Am. Soc. for Information Science and Technology, vol. 54, no. 7, pp. 638-649, 2003.

[8] T. Joachims, "Evaluating Retrieval Performance Using Clickthrough Data," Text Mining, J. Franke, G. Nakhaeizadeh, and I. Renz, eds., pp. 79-96, Physica/Springer Verlag, 2003.

[9] T. Joachims, "Optimizing Search Engines Using Clickthrough Data," Proc. Eighth ACM SIGKDD Int'l Conf. Knowledge Discovery and Data Mining (SIGKDD '02), pp. 133-142, 2002.

[10] T. Joachims, L. Granka, B. Pang, H. Hembrooke, and G. Gay, "Accurately Interpreting Clickthrough Data as Implicit Feedback," Proc. 28th Ann. Int'l ACM SIGIR Conf. Research and Development in Information Retrieval (SIGIR '05), pp. 154-161, 2005.

[11] R. Jones and K.L. Klinkner, "Beyond the Session Timeout: Automatic Hierarchical Segmentation of Search Topics in Query Logs," Proc. 17th ACM Conf. Information and Knowledge Management (CIKM '08), pp. 699-708, 2008.

[12] R. Jones, B. Rey, O. Madani, and W. Greiner, "Generating Query Substitutions," Proc. 15th Int'l Conf. World Wide Web (WWW '06), pp. 387-396, 2006.

[13] U. Lee, Z. Liu, and J. Cho, "Automatic Identification of User Goals in Web Search," Proc. 14th Int'l Conf. World Wide Web (WWW '05), pp. 391-400, 2005.

[14] X. Li, Y.-Y Wang, and A. Acero, "Learning Query Intent from Regularized Click Graphs," Proc. 31st Ann.

Int'l ACM SIGIR Conf. Research and Development in Information Retrieval (SIGIR '08), pp. 339-346, 2008.

[15] M. Pasca and B.-V Durme, "What You Seek Is what You Get: Extraction of Class Attributes from Query Logs," Proc. 20th Int'l Joint Conf. Artificial Intelligence (IJCAI '07), pp. 2832-2837, 2007.

[16] B. Poblete and B.-Y Ricardo, "Query-Sets: Using Implicit Feedback and Query Patterns to Organize Web Documents," Proc. 17th Int'l Conf. World Wide Web (WWW '08), pp. 41-50, 2008.