

SURVEY ON INFERRING USER SEARCH GOAL USING FEEDBACK SESSION

D.Kavitha, K.M.Subramanian, Dr.K.Venkatachalam

Abstract—Data Mining refers to extracting or “mining” knowledge from large amounts of data. It is also called as knowledge mining from data. Search engine is one of the most important applications in today’s internet. Users collect required information through the search engine in the internet. Analyzing user search goal is essential to provide best result for which the user looks for in the internet. Feedback sessions have been clustered to discover different user search goals for a query. Pseudo-documents are generated through feedback sessions for clustering. To understand the user search goals efficiently using Classified Average precision (CAP) algorithm.

Index terms—Data mining, user search goals, feedback sessions, pseudo-documents, classified average precision

I. INTRODUCTION

A Web is a collection of inter-related files on one or more Web servers. Web mining is the application of data mining technique it is used extract a knowledge from Web data. Web data is Web content data (text, image, record), Web structure data (hyperlinks, logs) and Web usage data (http logs, app server logs). In this paper we use the Web usage mining data. Discovery of meaningful patterns from data generated by client-server transactions on one or more Web localities. Analyzing and exploring regularities in Weblog records (consist of URL’s, time interval, click sequence and etc) for electronic commerce, enhance the quality and delivery of internet information services to the end user, and improve Web server system performance.

A Web server usually registers a log entry, or Weblog entry, for every access of a Web page. It includes the URL requested the IP address from which the request originated, and a timestamp. Based on the Weblog records. We have to construct the feedback session. Because Weblog data provide information about what kind of users will access what kind of Web pages. This session consist of RL’s and click sequence and it focus on user search goals. Only using a feedback session we do not understand the user search goals exactly.

Based on the feedback session construct the pseudo document for analyzing the accurate result. This pseudo document consist of keywords of URL’s in the feedback session. This is called as enriched URL’s. The enriched URL’s are clustered and form a pseudo document. Clustering is the process of grouping the data into classes or clusters, so that objects within a cluster have a high similarity in comparison to one another but are very dissimilar to object in other clusters. After constructing the pseudo document the Web search results are restructured based on the documents collection detail.

The rest of the paper is organized as follows: session 2 defines the detail about Web usage mining, session 3 defines the techniques like classification and prediction, and clustering, session 4 is detail about the

related work, and session 5 is about comparison method and session 6 about conclusion of this paper.

II. WEB USAGE MINING

Web usage mining is mines Weblog records to discover user access patterns of Web pages. In developing techniques for Web usage mining, we may consider the following

- It is encouraging and exciting to imagine the various potential applications of Weblog file analysis, it is important to know that the success of such applications depends on what and how much valid and reliable knowledge can be discovered from the large raw log data.
- The available URL, time, IP address, and Web page content information, a multidimensional view can be constructed on the Weblog database, and multidimensional OLAP analysis can be performed to find the top N users, top N accessed Web pages, most frequently accessed time periods, and so on which will help discover potential customers, users, markets, and others
- Data mining can be performed on Weblog records to find association patterns, sequential patterns, and trends of Web accessing.

III. CLASSIFICATION, PREDICTION AND CLUSTERING

The process of grouping a set of physical or abstract objects into classes of similar objects is called clustering. A cluster is a collection of data objects that are similar to one another within the same cluster and are dissimilar to the objects in other cluster. Dissimilarities are

assessed based on the attribute values describing the objects, often, distance measures are used. In this paper we use k-means clustering technique for constructing pseudo documents. K-means clustering is a centroid based technique. Classification and prediction are two forms of

data analysis that be used to extract models describing important data classes or to predict future data trends. Such analysis can help provide us with a better understanding of the data at large whereas classification predicts categorical labels, prediction models continuous-valued functions. It uses the preprocessing technique such as data cleaning, relevance analysis, data transformation and reduction. It provide the accuracy, scalability, robustness, speed and interpretability.

IV. RELATED WORK

A. Automatic Identification of User Goals in Web Search

Based on the Web query assigned by the user's analysis the goal, the goal identification is used to improve quality of search results. In existing system with use the manual query log investigation to identify the goals. In proposed system use automatic goal identification process. The human-subject study strongly indicates the automatic query goal identification. It can use two tasks like as past user click behavior and anchor link distribution for goal identification combining these two tasks can identify 90% goal accurately.

B. Query-Sets: Using Implicit Feedback and Query Patterns to Organize Web Documents

Document representation model (DRM) is based on the implicit user feedback. Implicit user feedback is mean that the feedback from weblog. Document representation model is obtained from search engine queries. The main objective of this DRM is to achieve the better results using non-supervised tasks such as clustering and labeling obtained from search engine queries. Users are motivated for document representation. Based on the clicked queries the term provide the better choice of feature

from the user's point of view. This model represent the frequency query patterns called as query set model. The query set model reduces the 90% the number of features needed for represent the set of documents, then improve 90% the quality of results.

C. Learn from Web Search Logs to Organize Search Results

Search results of the effective organization are critical to improve the utility of the search engine. Clustering the search results is the best way to organize the search results. Use the clustering of search results users finds the document quickly. There are two faults of this approach are: 1.The clusters do not depends on the interesting aspects of users. 2. The cluster labels are not informative, so that the identification of right clusters is difficult. The reasons are 1. Labels are not meaningful and 2. Labels are not informative. The solution of the faults in the proposed are: 1.Learning "interesting aspects" from web search logs and organizing search results. 2. Informative cluster labels are generated using query words used by the

D. Generating Query Substitutions

users. Evaluation of the method is based on commercial search engine log data. Compared with traditional method to this method to this method produce the better organization results and meaningful labels. The most common strategy of presenting search results is a simple

ranked list. Search engine logs record the activities of web users, which reflect the actual user's needs or interests when conducting web search. Search engine logs are separated by sessions. A session includes a single query and all the URLs that a user clicked after issuing the query.

E. Learning Query Intent from Regularized Click Graphs

Improve the query intent classifier using a click graphs, this method is critical for vertical and general purpose search services offered by user interface. In existing they use query classification for improving feature representation of queries. In proposed we focus on completely orthogonal approach for enriching feature representation. The main objective is to increasing the amounts of training data using semi-supervised learning with click graphs. Based on the click graph we understand the unlabeled queries from those of labeled ones. Moreover we regularize the learning with click graphs using content based classification to avoid the error labels. We define the effectiveness of our algorithms using two different application (product intent and job intent classification). Using this both applications we expands the training data and leading to improvements in classification performance. An additionally finding the large amount of training data based and classifiers using query words as features.

Query substitution generates the new query to replace the user's original query. This technique uses modification based on query substitution. The new queries

and the terms are closely related to the original queries and the terms. Query substitution is contrast with query expansion and query relaxation, the query expansion through pseudo-relevance feedback this is cost and lead to aimless process. The query relaxation through Boolean or TF-IDF retrieval, this reduces the specificity. Evaluation of query substitution is well performed the replacement of

F. Varying Approaches to Topical Web Query Classification

Web queries are classified based on the behaviors or some similarities. This classification of query improving retrieval effectiveness and efficiently. The query is used to retrieving a document before or after a query classification. We examine two previously unaddressed issues in query classification: 1. pre vs. post-retrieval classification, effectiveness and the effect of training explicitly from classified queries vs. bridging a classifier trained using a document taxonomy, 2. Bridging classifier maps the document taxonomy onto query classification problem and

it provide sufficient training data. We find that training classifier explicitly from manually classified queries to the

bridged classifier by 48% in F1 score. The pre-retrieval classifier is 11% worse than bridged classifier. It requires snippets from retrieved documents.

G. Context-Aware Query Suggestion by Mining Click-Through and Session Data

QS Plays an important role in improving the usability of search engine. In existing QS by mining query patterns from search logs, none of them are context aware.

In this paper the context in QS consist of two steps like 1. In offline process the learning step is used to address the data, queries are converted into concepts by a technique called clustering, a click through bipartite. Based on session data a sequence suffix tree is constructed for the QS model. 2. In online process the query suggestion is used to capture the user search results by mapping with the query sequence submitted by the user. This approach provides to the user in a context-aware manner. It is also called as Context-Aware Concept-Based Approach (CACBA)

H. Comparative Study

S.no	Titles	Technique/Methods	Advantage	Disadvantage
1.	Automatic Identification of User Goals in Web Search	User click behavior and Anchor link distribution	Using goal identification task to achieve 90% of accurate results	Potentially-biased dataset

2.	Query-Sets: Using Implicit Feedback and Query Patterns to Organize Web Documents	Non supervised tasks	Improve the quality as 90%	A broader comparison with online directory
3.	Learn from Web Search Logs to Organize Search Results	Commercial search engine log data and clustering	Better result organization and meaningful labels	Informative feedback information from user
4.	Generating Query Substitutions	Query pair algorithm	Increase coverage and effectiveness	Machine translation techniques
5.	Learning Query Intent from Regularized Click Graphs	Semi-supervised click graph	Improve classification performance	Impact of seed queries and faceted query classification
6.	Varying Approaches to Topical Web Query Classification	Pre vs. post retrieval classification	QC is outperforms bridging a document taxonomy as 48%	multiple approaches to improve performance
7.	Context-Aware QS by Mining Click-Through and Session Data	Offline model learning and online QS step, concept sequence suffix	Coverage and quality of suggestions	Larger coverage area

V. CONCLUSION

As the Web and its usage continues to grow, so grows the opportunity to analyze Web data and extract all manner of useful knowledge from it. The past five years have seen the emergence of Web mining as a rapidly growing area, due to the efforts of the research community as well as various organizations that are practicing it. In this paper, we propose a novel approach for user search goals using feedback session and pseudo document. First we construct a feedback session to analysis the user search goal from the Weblog record. It cannot provide the accurate

result. So that we can introduce the pseudo document to provide the accurate results. Based on the pseudo document we have to restructure the Web search results.

REFERENCES

- [1] Beitzel, S, 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.

- [2] Baeza-Yates. R, 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] Beeferman. D 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] Cao. H, 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.
- [5] Chen. H 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.
- [6] Huang C.K, 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.
- [7] Jones. R, B. Rey, O. Madani, and W. Greiner, "Generating Query Substitutions," Proc. 15th Int'l Conf. World Wide Web (WWW '06), pp. 387-396, 2006.
- [8] Joachims. T, 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.
- [9] Joachims. T, "Evaluating Retrieval Performance Using Clickthrough Data," Text Mining, J. Franke, G. Nakhaeizadeh, and I. Renz, eds., pp. 79-96, Physica/Springer Verlag, 2003.
- [10] Joachims. T, "Optimizing Search Engines Using Clickthrough Data," Proc. Eighth ACM SIGKDD Int'l Conf. Knowledge Discovery and Data Mining (SIGKDD '02), pp. 133-142, 2002.
- [11] Jones. R 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] Lee. U, 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.
- [13] Li. X, 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.
- [14] Poblete. B 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.
- [15] Pasca. M 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] Shen. D, J. Sun, Q. Yang, and Z. Chen, "Building Bridges for Web Query Classification," Proc. 29th Ann. Int'l ACM SIGIR Conf. Research and Development in Information Retrieval (SIGIR '06), pp. 131-138, 2006.

- [17] Wen J-R, J.-Y Nie, and H.-J Zhang, “Clustering User Queries of a Search Engine,” Proc. Tenth Int’l Conf. World Wide Web (WWW ’01), pp. 162-168, 2001.
- [18] Wang. X and C.-X Zhai, “Learn from Web Search Logs to Organize Search Results,” Proc. 30th Ann.

Int’l ACM SIGIR Conf. Research and Development in Information Retrieval (SIGIR ’07), pp. 87-94, 2007.

- [19] Zeng. H, Q.-C He, Z. Chen, W.-Y Ma, and J. Ma, “Learning to Cluster Web Search Results,” Proc. 27th Ann. Int’l ACM SIGIR Conf. Research and Development in Information Retrieval (SIGIR ’04), pp. 210-217, 2004.

University, Gundiy and Presently he is working as Professor in Department of Electronics and Communication Engineering, Velalar College of Engineering and Technology, Erode. His Area of Specialization includes Wireless Mobile Communication, Optical Networks, Wireless Embedded Networks and Low Power VLSI Design.



Kavitha D received B.E degree in computer science and Engineering from Nandha College of Technology, Erode. She is currently pursuing her M.E degree in computer science and Engineering in Erode Sengunthar Engineering College. She has published 2 papers in national conferences. Her Research interest includes Data Mining and Database Management System.



Subramanian K.M received his B.Sc. degree in Computer Science from Cheran Arts & Science College, Kangayam (Bharathiar University), Tamilnadu. The M.C.A degree from Bharathidasan University, Tiruchirappalli, and M.E. Degree in Computer Science & Engineering from Mahendra Engineering College, Salem. His research interests include Data Mining, Data Warehousing, Web Mining and Neural Networks. At present he is working as a Assistant Professor – Selection Grade-II, Department of Computer Science & Engineering, Erode Sengunthar Engineering College, and pursuing his Ph.D., degree in Anna University, Chennai, India. He has published 3 international journals and 2 national journals. He is a ISTE Life member and CSI Institutional member.



Dr. Venkatachalam K received B.E degree in Electronics and Communication Engineering from Bharathiar University, Coimbatore, M.Tech degree in Electronic and Communication System- First class with Distinction and Universty 2nd rank from Pondicherry University and Ph.D degree in Information Communication Engineering from Anna