

GLOBAL JOURNAL OF ENGINEERING SCIENCE AND RESEARCHES MINING SPECIFIC CONTENT SEARCH USING CLUSTERING ALGORITHM

Y.Padma^{*1}, P.RaviPrakash² & D.Kavitha³

^{*1&2}Asst. Professor, Dept. of IT, PVPSIT, Vijayawada

³Sr. Asst. Professor, Dept. of IT, PVPSIT, Vijayawada

ABSTRACT

Recent challenges in information retrieval are related to information in social networks and rich media content. In those cases, the content is associated with multilingual, user generated aspects and content, scalability, robustness and resilience to errors. The graphical model learns is likely a candidate term is to be a facet term as well as how likely two terms are to be grouped together in a query facet, and captures the dependencies between the two factors. Proposed system enhanced the previous work to avoid duplication of similar site by page parsing and comparison of page content. We propose a clustering algorithm is effectively leverage the two phenomena to automatically mine the major subtopics of queries where each subtopic is represented by a cluster containing a number of URLs and keywords. Moreover fast and efficient indexing and searching services are needed, in order to scale digital content distribution and video on demand, where huge amount of queries and content related tasks are performed by users. To estimate the size of a hidden database, on intuitive idea is to performtopple sampling.

Keywords- *Query Facet, Semantic Class Extraction, Web crawling, indexing, QD miner, Query Subtopics, Clustering, Search Result Clustering.*

I. INTRODUCTION

Understanding the search intent of users is essential for satisfying a user's search needs it best represent query intent is still an ongoing research problem. One consensus in the researchers is that the intents of queries is characterized along multiple dimensions [1]. The intents of a query is represented by its search goals, such as informational, navigational, and transactional [2]. We develop supervised method based on a graphical model for query facet extraction. The graphical model learns a term should be selected and how likely it is that two terms should be grouped together in a query facet [3]. The model captures the dependencies between the two factors. The concept of designated query to form an injective mapping from tuples to queries supported by the web interface [4]. To produce unbiased aggregate estimations over the hidden databases with checkbox interfaces, develop the data structure of left-deep-tree and define the concept of designated query to form an injective mapping from tuples to queries supported by the web interface [5]. Typically, there is a need of tools for metadata extraction, schemas and metadata mapping rules and tools, multilingual metadata and content translation and certification. Information retrieval (IR) systems are required to give coherent answers with respect to typos or inflexions, and must be efficient enough while sorting huge result lists [6]. Context Similarity Model, in which we model the filtered similarity between each pair of product. Link classification suggests user interest content mining in various aspects like shopping, education, searching [7].

II. RELATED WORK

The topics subtopics of a query are represented by a number of queries or URLs. Topics is usually more coarse-grained and cover multiple queries, while subtopics are more fine-grained and associated with a specific query [8]. Mining topics from search log data has been intensively studied. Click-through bipartite graph data is used for clustering queries and URLs. Specifically queries which share the same clicked URLs are considered similar. Methods for performing the task have been proposed [9]. Proposed conducting clustering on a click-through bipartite graph and viewing the obtained clusters as topics covering multiple queries author represents query facets to understand user interest for search in diversification [10]. Investigated model generates subtopics based on query factors and proposed faceted diversification approaches. Each facet contains a group of words or phrases extracted

from search results [11]. Faceted search is a technique for accessing information organized according to a faceted classification system, acceding users to digest, analyze and navigate through multidimensional data. It is new used in e-commerce and digital libraries [6]. Faceted search is similar to query facet extraction in that both of them use sets of coordinate terms to represent different facets of a query [12]. The Deep Web is also believed to be the biggest source of structured data on the Web and hence accessing its contents has been a long standing challenge in the data management community .Over the past few years they have built a system that exposed content from the Deep Web to web-search users of Google.com [13].

III. SYSTEM ARCHITECTURE

Query facet mining from huge searchable data is cumbersome task. This system modify facet mining task with the help of natural language processing for HTML form data.Methods for utilizing a user click behavior in different searches have been developed [15]. The exploitation of the prefix and suffix relationship in queries is considered in the previous work. In our work not only use the prefix and suffix relationship between queries, but also the clicked URLs of the queries and our goal is to conduct query subtopic mining [14].

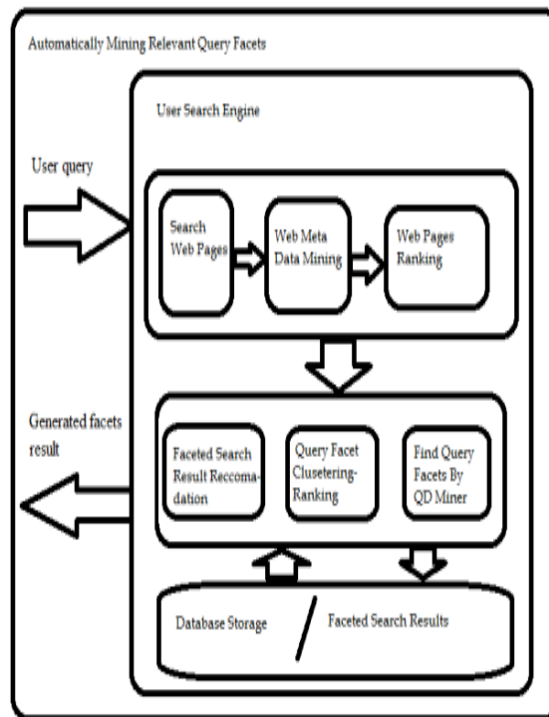


Fig.1: System architecture of automatically Query facet mining

IV. PROPOSED APPROACH

Search result clustering attempts to cluster the search results according to semantic classes, topics. The clicked URLs after searching with the original query and the expanded queries model to represent the same subtopic [16]. This is phenomenon of subtopic clarification by additional keyword.Furthermore number of hidden databases does not publicize their total sizes, while such information is useful to the general public as an economic indicator for monitoring product growth [17].

A. Graphical model

We define all the variables in our graphical model. Let $Y = \{y_i\}$, where $y_i = 1\{t_i \in TF\}$ is a label indicating whether a list item t_i is a facet term. Here $1\{\cdot\}$ is an indicator function which takes on a value of 1 if its argument is true, and 0 otherwise. $p_{i,j}$ denotes the list items pair (t_i, t_j) , and $PL = \{p_{i,j} | p_{i,j} = (t_i, t_j), t_i, t_j \in TL, t_i \neq t_j\}$ denotes all the items pairs in TL. Let $Z = \{z_{i,j}\}$, where $z_{i,j} = 1\{\exists F \in F, t_i, t_j \in F \wedge t_j \in F\}$ is a label indicates [19] whether the corresponding item pair $p_{i,j}$ should be grouped together in a query facet. The vertices in our graphical model are $V = TL \cup PL \cup Y \cup Z$. Note that the list items TL, and item pairs PL are always observed.

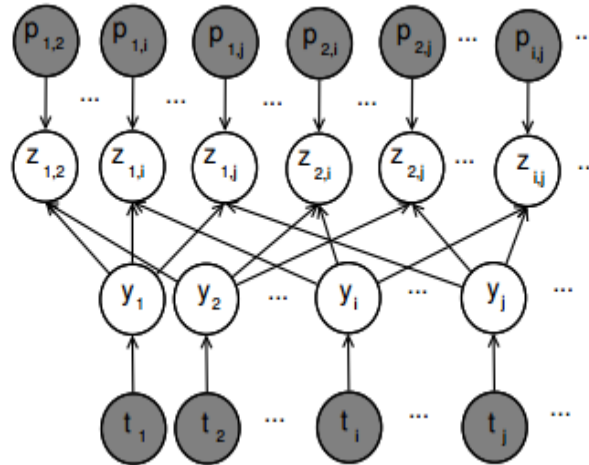


Figure 2: A graphical model for candidate list data

The algorithm is summarized it processes the facet terms in decreasing order of $P(t)$. For each facet term remaining in the pool, it builds a cluster by iteratively including the facet term that is closest to the cluster, until the diameter of the cluster surpasses the threshold d_{max} [18].

Algorithm: WQT for clustering facet term used in QF-I

Input: TF $P(t)$, $d_f(F, t)$, $dia(F)$, d_{max}

Output: $F = \{F\}$

- 1: $T_{pool} \leftarrow F$
- 2: repeat
- 3: $t \leftarrow \arg \max_{t \in T_{pool}} P(t)$
- 4: $F \leftarrow \{t\}$
- 5: iteratively include facet term $t_0 \in T_{pool}$ that is closest to F , according to $d_f(F, t_0)$, until the diameter of the cluster, $dia(F)$, surpasses the threshold d_{max} .
- 6: $F \leftarrow F \cup \{F\}$, $T_{pool} \leftarrow T_{pool} - F$
- 7: until T_{pool} is empty
- 8: return F

V. CLUSTERING METHOD

Clustering method to mine subtopics of queries leveraging the two phenomena and search log data. We build an index to store all the queries and their clicked URLs. False expanded queries are then pruned from the index. In the clustering stage, the URLs associated with a query and its expanded queries are grouped into clusters, each representing one subtopic [19].

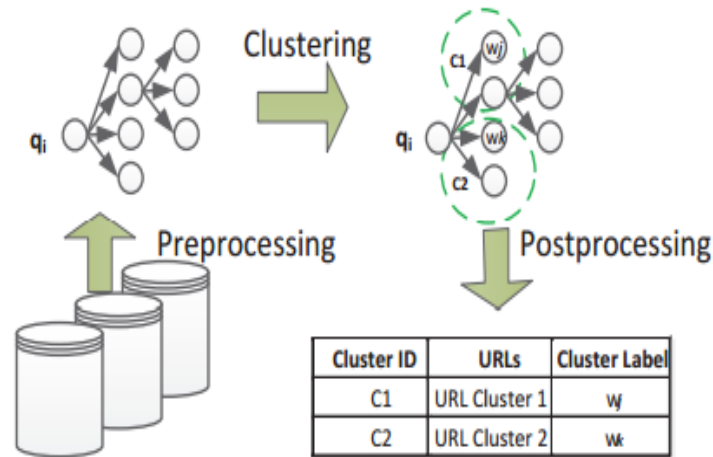


Figure 3: The flow of clustering method

VI. INDEXING

We first index all the queries in an index consisting of a prefix tree and a suffix tree to facilitate efficient clustering. We only consider queries in three forms ('Q', 'Q + W' and 'W + Q'), We then segment queries and index them. In the prefix tree, query 'Q' and its expanded queries 'Q+W' are indexed in a father node and child nodes respectively [20]. With the prefix tree the suffix tree query 'Q' and its expanded queries 'W+Q' are indexed as a father node and child nodes respectively

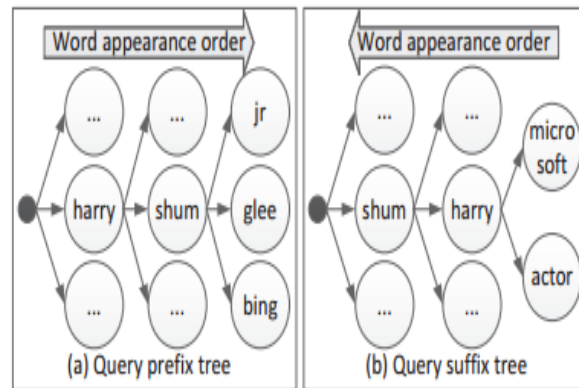


Figure 4: The data structures to index search logs

VII. CLUSTERING

We conduct clustering on the clicked URLs of each query and its expanded queries. Since all the queries are indexed in the trees, the clustering is performed locally and recursively on the trees. The clustering of clicked URLs is guided by the two phenomena. After clustering, each group of clustered URLs is taken as one subtopic of the query in the father node [21].

Algorithm

We employ an agglomerative clustering algorithm to perform clustering. The algorithm has the advantage of ease of implementation. One can also consider employing other clustering algorithms. The specific algorithm is as follows:

Step 1: Select one URL and create a new cluster containing the URL.

Step 2: Select the next URL u_i , and make a similarity comparison between the URL and all the URLs in the existing clusters. If the similarity between URL u_i and URL u_j in one of the clusters is larger than threshold θ , then move u_i into the cluster [22].

Step 3: Finish when all the URLs are processed.

VIII. SEARCHING

The goal of the searching service is access the users to easily locate and sort each type of content in the ECLAP portal, and to refine their queries for a more detailed result filtering, through a fast search interface, robust with respect to mistyping; high granularity of data has to be offered to the users.

Boosting of Terms is configurable on the portal. This allowed us to tune and stress the importance of certain metadata. Boosting and weighting of metadata are better tuned when the portal is more populated with significant contents. Each field of the ECLAP document structure is boosted with its predefined value at query time [23]

Faceted search is activated on the results of both simple frontal search and advanced search. Each faceted term is indexed un-tokenized in the ECLAP index, to accomplish a faceting count based on the whole facet. Drupe service module before rendering. The user can select or remove any facet in any order to refine the search a search filter, and performs again the search query with or without it. Relevant facets include:

- A. DC: resource category, format, type, classification, creator, content language, etc.
- B. Technical: duration, video quality, device, publisher source metadata language and upload time
- C. Group, taxonomy: genre, historical period, performing arts, coded subject

These facets can be subject to change. For instance, locations and dates, different for each historical period, can be added [24].

IX. SEARCH RESULTS

Search results are listed by relevance in descending order this means that the first document is the most relevant with respect to the query. The relevance is based on the occurrence of the query term in the indexed document fields a higher number of term's occurrences give a higher score for the document. Each result item is presented with a thumbnail, relevant metadata rating, relevance score and number of accesses; data is presented in the same language chosen by the user among the available portal localizations.

users	# Full Text Queries	# of Faceted Queries	# Last Posted Contents	# Featured Contents	# Popular Contents
simple registered	323	24	4	22	17
Registered as partners	1094	21	27	19	9
anonymous	2634	147	234	302	213
Total	4051	192	265	343	239
Clicks after query	1564	200	318	2799	231

Fig no 5. QUERIES / CONTENT LISTS

It can be noted that after a query on the portal, the 92.65% search results clicks were performed in the first page (first ten results). 42.27% of clicks on search results have been performed to the first proposed result. The second has received only the 14% clicks.

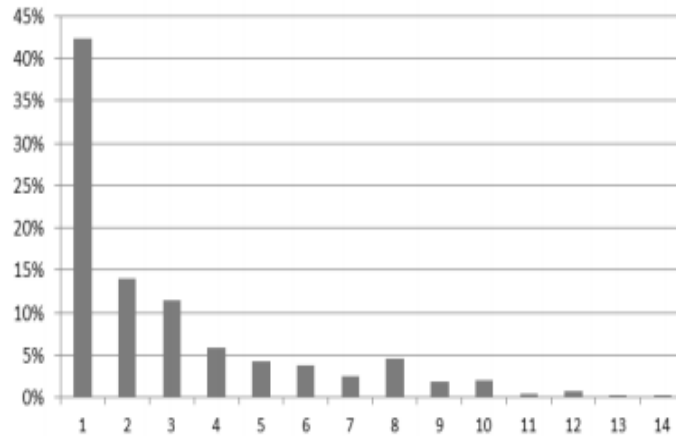


Figure 6. Clicks order distribution (first page results and a part fo the second)

X. CONCLUSIONS

We developed a supervised method based on a graphical model to recognize query facets from the noisy facet candidate lists extracted from the top ranked search results. The scope for probing the hidden databases since query probing techniques is widely used in the hidden database. Proposed mining access fine grained facets from search result for client search query relevant URLs is gathered by applying reverse search algorithm and indexing the accessible report by naive Bayes classifiers. We have developed a clustering algorithm and effectively and efficiently mine query subtopics on the basis of the two phenomena. We have evaluated the effectiveness of the proposed models. Finally, our method is employed search log data, which is also a drawback for most log mining algorithms to apply the approach in tail queries is also an issue we need to consider.

REFERENCES

- [1] O. Ben-Yitzhak, N. Golbandi, N. Har'El, R. Lempel, A. Neumann, S. Ofek-Koifman, D. Sheinwald, E. Shekita, B. Sznajder, and S. Yogev, "Beyond basic faceted search," in *Proceedings of WSDM '08*, 2008.
- [2] A. Z. Broder. A taxonomy of web search. *Sigir Forum*, 36:3–10, 2002.
- [3] N. Bansal, A. Blum, and S. Chawla. Correlation clustering. In *MACHINE LEARNING*, pages 238–247, 2002.
- [4] C. Carpineto, S. Osiński, G. Romano, and D. Weiss. A survey of web clustering engines. *ACM Comput. Surv.*, 41(3):17:1–17:38, July 2009.
- [5] S. Jiantao, "Analyzing the network friendliness of mobile applications," Huawei, Shenzhen, China, Tech. Rep. M3-001034414-20120731-C-2.0, Jul. 2012
- [6] G. Gorbil and E. Gelenbe, "Resilience and security of opportunistic communications for emergency evacuation," in *Proc. 7th ACM Workshop Perform. Monitor. Meas. Heterogeneous Wireless Wired Netw. (PM2HW2N)*, Oct. 2012, pp. 115–124.
- [7] T. Taleb and A. Kunz, "Machine type communications in 3GPP networks: Potential, challenges, and solutions," *IEEE Commun. Mag.*, vol. 50, no. 3, pp. 178–184, Mar. 2012.
- [8] A. Z. Broder, M. Fontoura, E. Gabrilovich, A. Joshi, V. Josifovski, and T. Zhang. Robust classification of rare queries using web knowledge. In *Proceedings of SIGIR '07*, pages 231–238, 2007.
- [9] M. Burt and C. L. Liew. Searching with clustering: An investigation into the effects on users' search experience and satisfaction.
- [10] Mahinthan Chandramohan and HeeBengKuan Tan, "Detection of Mobile Malware in the Wild," 0018-9162/12/\$31.00 © 2012 IEEE Published by the IEEE Computer Society SEPTEMBER 2012.
- [11] Erol Gelenbe, Gokce Gorbil, Dimitrios Tzovaras†, Steffen Liebergeld, "Security for Smart Mobile Networks: The NEMESYS Approach," *Mobile Telecommunications S.A.*, 15124 Maroussi, Greece
- [12] Homefinder, Home finder page [Online]. Available: <http://www.homefinder.com/search>, 2013.

- [13]A. Dasgupta, X. Jin, B. Jewell, N. Zhang, and G. Das, *Unbiased estimation of size and other aggregates over hidden web databases*, in *Proc. Int. Conf. Manage. Data*, 2010, pp. 855–866
- [14]C. L. A. Clarke, N. Craswell, and I. Soboroff. *Overview of the trec 2009 web track*. In *Proceedings of TREC'09*, pages 1–9, 2009.
- [15]N. Craswell and M. Szummer. *Random walks on the click graph*. In *Proceedings of SIGIR'07*, pages 239–246, 2007
- [16]H. Deng, I. King, and M. Lyu. *Entropy-biased models for query representation on the click graph*. In *Proceedings of SIGIR'09*, pages 339–346. ACM, 2009.
- [17]P. Ferragina and A. Gulli. *A personalized search engine based on web-snippet hierarchical clustering*. *Software: Practice and Experience*, 38(2):189–225, 2008.
- [18]M. Pasca and E. Alfonseca. *Web-derived resources for web information retrieval: from conceptual hierarchies to attribute hierarchies*. In *Proceedings of SIGIR '09*, pages 596–603, 2009.
- [19]P. Pantel, E. Crestan, A. Borkovsky, A.-M. Popescu, and V. Vyas. *Web-scale distributional similarity and entity set expansion*. In *Proceedings of EMNLP '09*, pages 938–947, 2009.
- [20]T. Joachims. *Optimizing search engines using clickthrough data*. In *Proceedings of KDD'02*, pages 133–142, 2002.
- [21]R. Jones and K. Klinkner. *Beyond the session timeout: Automatic hierarchical segmentation of search topics in query logs*. In *Proceedings of CIKM'08*, pages 699–708, 2008.
- [22]M. Kan and H. Thi. *Fast webpage classification using url features*. In *Proceedings of the 14th ACM international conference on Information and knowledge management*, pages 325–326. ACM, 2005
- [23]G. M. Sacco, “*Research Results in Dynamic Taxonomy and Faceted Search Systems*”, *18th International Workshop on Database and Expert Systems Applications*, 2007.
- [24]D. Tunkelang, “*Faceted Search*”, (*Synthesis Lectures on Information Concepts, Retrieval, and Services*), Morgan and Claypool Publishers, 2009