# QDminer: Finding Facet Query by Modelling Fine Grained Similarities

Rupali Vasant Ubale
ME Student
Department of Computer Engineering
Dr. D. Y. Patil College of Engineering, Pune, India

**Abstract:**
Query facets offer interesting as well as beneficial knowledge about any query and hence it can be used to increase search practices in numerous ways. Thus in this paper, the issues regarding the finding of query facets are summaries. Query facets are different groups of words or else phrases which explain and abstract the content enclosed by a query. From the study, it is observed that the correlated aspects of a query are typically granted and also repeated in the top-k retrieved documents of the query in the form of document lists, and by aggregating these lists the query facets can be mined out. In this paper an efficient solution is projected which stated to as QDMiner, to automatically mine query facets by extracting then grouping common lists from HTML tags, free text and duplicate sections within topmost search results. Furthermore, the difficulty in the duplication of the list is examined and better query facets are found that can be mined by demonstrating fine-grained relationships among lists and fining the duplicated lists. In addition, the collaborative filtering techniques are used for recommendation of top-k results of user interest. In this recommendation process, ICHM and UCHM techniques are used to predict results according to user interest through matrix generation.

**Index Terms:** Query facet, Faceted search, Summarization, collaborative filtering.

## I. INTRODUCTION

### A. Query Facets

A query facet is a set of items like words or various phrases which define and summarize significant characteristic of a query. One single query may have numerous facets that describe the information regarding the query from different viewpoints. The query"visit Mumbai" has a query facet about popular hotels in Mumbai (Marine Drive, hotel Taj, Gate way of India, . . .) and a facet on travel associated areas (attractions, supermarket run, dining, . . .). Query facets offer interesting as well as beneficial knowledge about any query and hence it can be used to increase search practices in numerous ways. The users can simplify their exact intent by choosing facet things. Then search results might be limited to the documents that are significant to the things. A user possibly will drill down to ladies watches if he is watching for a gift for his wife. These numerous groups of query facets are in specific suitable for unclear or confusing queries, such as "apple". We might display the products of Apple Inc. in any facet as well as various types of the apple fruit in another facet. Another, query facets could deliver direct information or immediate answers that users are looking for. Such as, for the query "big boss season 6", all episode designations are displayed in one facet and leading actors are displayed in another. In this circumstance, displaying query facets might save browsing time. Third, query facets might be used to increase the variety of the ten blue links. We may re-rank search results to avoid displaying the pages which are close to replicated data in query facets on the top. Furthermore, query facets comprise structured information enclosed by the query, and therefore query facets can be used in new fields besides customary web searches, for example, semantic search or else entity search.

### B. Collaborative Filtering

Collaborative filtering (CF) is a standard technology for recommender schemes. These methods are classified into two parts one is user-based CF and another is item-based CF. The main goal of user-based CF approach is to search out a group of users who have similar favor forms to a known user (i.e."neighbors" of the user) and recommend those things to the user that different users in the same set, while the item-based CF approach goals to offer a user with the recommendation on an item supported the opposite things with high correlations (i.e. "neighbors" of the item). In entirely collaborative filtering strategies, it is a major step to finding user's (or item's) neighbors, i.e., a set of comparable users (or items). At present, almost complete CF strategies measure user's similarity (or thing's similarity) supported co-rated items of users (or collective users of items). Though these recommendation strategies are widely used in E-Commerce, a number of inadequacies are identified. Recently we are typically overwhelmed by the large volume of information accessible on the net, and in this environment, we should build choices relating to the consumption of data. In our everyday survives, critics do much of our data filtering. For example, check rank lists for blockbusters and pay attention to movie critics. Collaborative filtering scheme overwhelms some restrictions of content-based filtering. The method can recommend items (such as music, books etc.) to users as well as recommendations are built on the ratings given to the items, as an alternative of the contents of the items, which can increase the quality of recommendations. Though collaborative filtering has been effectively used in together research and exercise, there still continue nearly challenges for it as an effectual data filtering. From the earlier study, it is perceived that significant sections of data about a query are typically accessible

in list formats and frequently used numerous times between top-k retrieved documents. Therefore frequent lists aggregating inside the top-k search results are planned to mine query facets as well as implement a method called as QDMiner. More precisely, QDMiner retrieves lists from HTML tags, text as well as replicate regions enclosed in the top-k search results, combines them into clusters depends on the items they enclose, then orders the clusters as well as items based on in what manner the lists and items seem in the top-k results. The scheme comprises two representations, one is the Unique Website Model, and another is the Context Similarity Model, to order or rank query facets. Furthermore, to recommend user interested result, a collaborative filtering technique is used. As for a collaborative recommendation, there are two ways to calculate the similarity for group recommendation: Item based and user-based. The next sections of the paper are organized as follows: Section II gives the important literature survey. Section III addresses proposed system. Section IV introduces the process in mathematical manner of the proposed system. Section V describes assumptions expected results. Section VI accomplishes the paper.

## II. REVIEW OF LITERATURE

In the literature review, we are going to discuss topical methods over the collaborative filtering and query facet search. In [1] L. Bing et al. suggest a graphical model to give score queries. The suggested model feats a latent topic space, which is automatically resulting from the log of query, to identify semantic dependence of terms in a query as well as dependency between topics. The graphical model correspondingly captures the context of term in the history query through skip-bigram in addition to n-gram language models. W. Kong et al. [2] challenge the heterogeneous environment of the web suggests to use automatic query-dependent facet generation, which creates facets for a query as an alternative of the entire corpus. To integrate feedback of user on these query facets into document ordering, they investigate together Boolean filtering as well as soft ranking models. I. Szpektor et al. [3] recommend a technique to extend the influence of query assistance methods as well as specific query recommendation to long-tail queries by thinking about rules among query patterns instead of individual query evolutions, as presently done in graph models of query-flow. X. Xue and W. B. Croft [4] projected a framework that represents reformulation as a distribution of queries, where each query is a variant of the actual query. This methodology deliberates a query as a simple unit and may capture significant dependencies among words as well as phrases in the query. Preceding reformulation models are different cases of the projected framework by creating particular assumptions. L. Liet al. [5] projected the three-phase framework designed for personalized query recommendations. The primary phase is the training of queries and their significant search results returned by a search engine, which creates a historic queryURL bipartite group. The next phase is the finding of related queries by retrieving a query affinity graph from the bipartite graph, rather than directly working on the original bipartite graph using biclique-based methodology or graph clustering. The third phase is to rank or order the similar queries. For this phase they create a rank technique for ordering the associated queries based on the merging distances of a hierarchical agglomerative clustering (HAC). W. Kong [6] improves a supervised method built on a graphical model to identify and recognize query facets from the noisy candidates found. The graphical model studies in what manner possibly a candidate term is to be a facet term along with how probable two terms are to be gathered together in a query facet also captures the dependencies among the two factors. They suggest two procedures for approximate implication on the graphical model temporarily exact inference is inflexible. Qing Li et al. [7] applied a clustering method to assimilate the contents of things into the framework of the item-based collaborative filtering. The group rating data that is achieved from the clustering outcome delivers a mode to present content data into a collaborative recommendation. Szpektor et al. [8] suggest a method to encompassthe reach of query support methods and in specific query recommendation to long-tail queries by reasoning nearby rules among query templates instead of individual query transitions, as presently done in graph models of query-flow.

I. Pound et al. [9] presented the user faceted-search performance using the connection of web query logs with present structured information. Meanwhile web queries are expressed as free-text queries; a challenge in this method is the inherent ambiguity in mapping keywords to the dissimilar probable attributes of assumed entity type. They present a solution that produces user partialities on attributes as well as values, employing dissimilar disambiguation methods ranging from humble keyword matching to additional sophisticated probabilistic models.

M. Diao et al. [10] apply the ideas of faceted search in addition browsing to the SpokenWeb search issue. They use the ideas of facets to index the metadata related to the audio content. Authors deliver a mechanism to order the facets created on the search results. They develop a collaborative query interface that allows browsing of search outcomes over the top ranked facets. K. Balog et al. [11] deliberate the task of entity search as well as study to which extent state-of-art information retrieval (IR) and semantic web (SW) skills are accomplished of answering data requirements that focus on entities. They similarly explore the possibility of merging IR with SW technologies to increase the end-to-end presentation on explicit entity search task. M. Bron et al [12] examined the presentation of a model that individual uses co-occurrence statistics. Though it recognizes a set of associated entities, it fails to order them efficiently. Two types of error arise: (1) entities of the incorrect type contaminate the ranking then (2) though some-how related to the basic entity, some extracted entities do not involve in the correct relation to it. To address error (1), they enhance type filtering based on group information accessible in Wikipedia. To precise for (2), they improve contextual data, characterized as language models resulting from documents in which source as well as destination entities co-occur. To finalize the pipeline, they find homepages of top-ranked entities by merging a language modeling method with heuristics established on Wikipedia's outer links. C. Li et al. [13] suggests Faceted media, a faceted recovery method for data discovery in addition exploration in Wikipedia. Assumed the group of Wikipedia articles subsequent from a keyword query, Facetedmedia creates a faceted interface for navigation of the result articles. Associated with other faceted retrieval methods, Facetedmedia is completely automatic as well as dynamic in together facet generation as well as hierarchy construction, and the facets are created on the rich semantic documents from Wikipedia. A. Herdagdelen et al. [14] offered method to query reformulation which associates syntactic as well as semantic data by means of generalized Levenshtein distance algorithms where

the replacement process costs are grounded on probabilistic term rewrite functions. They examine unsupervised, compact and effectual models, as well as deliver empirical evidence of their efficiency. They additional discover a query reformulation generative model and supervised grouping approaches providing better performance at variable computational costs. J. Huang and E. N. Efthimiadis [15] learning user's reformulation approaches in the context of the AOL query logs. They generate the taxonomy of query refinement approaches and construct a high precision rule-based classifier to identify separately type of reformulation. The efficiency of reformulations is dignified using user click activities. S. Gholamrezazadeh et al. [16] offerings the taxonomy of summarization schemes and describes the most significant criteria for a summary that may be produced by a scheme. Moreover, dissimilar approaches of text summarization, besides key steps for summarization procedure are deliberated. Likewise, go over core criteria for calculating a text summarization. H. Zhang et al. [17] studies the employment of topic models to construct semantic classes, taking as the basis data a collection of raw semantic classes (RASCs), which were mined by applying prescribed designs to web pages. The main necessity and challenge is to deal with multi-membership: An item could belong to numerous semantic classes, and need to determine many conceivable the dissimilar semantic classes the item belongs to. They treat RASCs by way of "documents", items by way of "words" and the last semantic classes by way of "topics" to accept topic models. O. Ben-Yitzhak [18] extends faceted search to support comfortable data detection responsibilities over more complex information models. Their primary extension enhances flexible, dynamic business intelligence combinations to the faceted presentation, allowing users to increase insight into their information that is far richer than impartial knowing the numbers of documents going to respectively facet. They understand this ability as a step toward bringing OLAP abilities, conventionally supported by databases completed relational information, to the domain of free-text queries over metadata-rich content. Their next extension displays how one can capably extend a faceted search engine to provision associated facets additional complex data model in which the values related with a document through multiple facets are not independent. W. Dakka and P. G. Ipeirotis [19] detect that facet terms rarely perform in text documents, viewing that they require exterior resources to recognize useful facet terms. For this, they first identify significant phrases in respective document. Then, they develop respective phrase with "context" phrases by means of external assets, for example WordNet and Wikipedia, producing facet terms to perform the extended database. Lastly, they associate the term deliveries in the original database then the extended database to recognize the terms that might be used to build browsing facets. S. Riezler et al. [20] apply pairs of user queries as well as snippets of user clicked results to train a model of machine translation to associate the "lexical gap" among query and document space. They show that the combination of a query-to-snippet translation model through a huge n-gram language model trained on queries accomplishes developed relative query extension associated to a method based on term correlations.

## III. PROPOSED SYSTEM

In QDMiner, for a query $q$, the top-k results are retrieved from a search engine then fetch complete documents to form a set $R$ as input. After that, query facets are excavated by four methods:

- Extraction of list and context: Lists and their context are mined from every document in set $R$. "men's watches, kid's watches, women's watches, luxury watches," is a sample list mined.
- List weighting: Completely extracted lists are weighted, then therefore some insignificant or noisy lists infrequently occur on a page, for example the price list "290.99, 340.99, 490.99...", that may be allocated by low weights.
- List clustering: Related lists are clustered collected to comprise a facet. Such as, different lists near watch gender categories are grouped since they share the similar items "men's" as well as "women's".
- Facet and item ranking: Further Facets as well as their items are assessed and ranked. Such as, the facet on brands is ordered higher than the facet on colors centered on in what way frequent the facets occur and how appropriate the associated documents are. Inside the query facet on gender types, "men's" then "women's" are placed higher than "unisex" and 'kids" built on how common the items seem, and their rank in the original lists.

This paper also suggests the technique that presents the contents of items into the item-based collaborative filtering to increase its prediction excellence and resolve the cold start difficulty. The technique is called as ICHM (Item-based Clustering Hybrid Method in which the item data and user ratings are combined to compute the item-item resemblance. Clustering method not simply can be applied to item-based collaborative recommenders but moreover may be applied to user-based collaborative recommenders. The technique is called as UCHM (User-based Clustering Hybrid Method which is based on the characteristics of user profiles as well as clustering consequence is preserved as items. Though, in ICHM, clustering is based on the characteristics of items and clustering result is preserved as users.

### A. Proposed Architecture Diagram
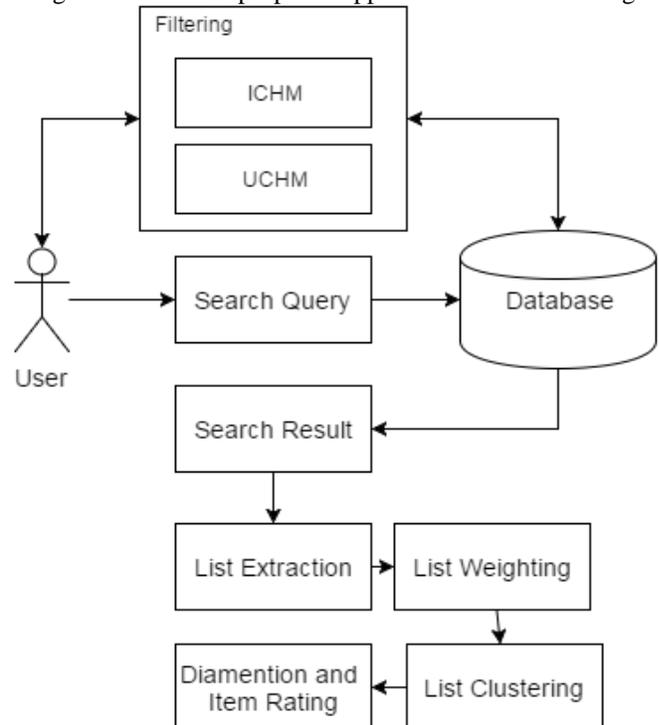Figure 1 shows the proposed approach for emotion recognition.



**Figure. 1. Architecture of Proposed System**

## IV. PROPOSED SYSTEM PROCESSING

### A. Preprocessing

**Input:**
D-set of documents $D = d_1, d_2, ..., d_n$ and $L_d$-set of lists $L_d = \{l^{0}\}$ extracted from the HTML content of $d$.

**Process:**

**1) List Weighting**

a) Compute document matching weight as:

$$S_{DOC} = \sum_{d \in R} \left( s_d^m \cdot s_d^r \right)$$

where, $s_{md}.s_{rd}$ is the supporting score by each result $d$,

$$s_d^m = \frac{N_{l,d}}{|l|}$$, where, $N_{l,d}$ is the number of items which appear both in list $l$ and document $d$

$$s_d^r = \frac{1}{rank_d}$$

$d$ , where $rank_d$ is the rank of document

b) Compute average invert document frequency (IDF) of items:

$$S_{IDF} - \frac{1}{|l|} \cdot \sum_{e \in l} idf_e$$,

where, $idf_e = log \frac{N - N_e + 0.5}{N_e + 0.5}$ , Where $N_e$ is the total number of documents that contain item $e$ in the corpus and $N$ is the total number of documents.

c) Evaluate the importance of a list $l$ as:
$Sl = SDOC.SIDF$

**2) List Clustering**

Use the complete linkage distance to compute the distance between two clusters of lists $l_1$, $l_2$. $dc(c1, c2) = maxl1 \in c1, l2 \in c2 dl(l1, l2)$

$$d_l(l_1, l_2) = 1 - \frac{|l_1 \cap l_2|}{min\{|l_1|, |l_2|\}}$$

Where,

**3) Facet Ranking**

a) Unique Website Model

Let $Cc = Sites(c)$ and recall that $Sites(c)$ is the set of unique websites containing lists in $c$. Then we have:
$Sc = Ps \in Sites(c) maxl \in c, l \in sSl$

b) Context Similarity Model

Similarity between two lists $l_1$ and $l_2$ is then calculated based on Hamming Distance $dist(l_1, l_2)$ between the fingerprints of their context:

$$Dup_L(l_1, l_2) = 1 - \frac{dist(l_1, l_2)}{LS}$$

**4) Item Ranking**

a) Calculate the weight of an item $e$ within a facet $c$ as:

$$S_{e|c} = \sum s \in C(c) w(c, e, C) = \sum_{G \in C(c)} \frac{1}{AvgRank_{c,e,G}}$$

Where $w(c, e, C)$ the weight contributed by a group of lists $G$, and $AvgRank_{c,e,G}$ is the average rank of item $e$ within all lists extracted from group $G$.

b) Suppose $L(c, e, G)$ is the set of all lists in $c$ and $G (G \subseteq c)$ that contain item $e$, we have

$$AvgRank_{c,e,G} = \frac{1}{|L(c, e, G)|} \sum_{l \in L(c,e,G)} rank_{e|l}$$

And $w(c, e, G)$ gets the highest score 1.0 when the item $e$ is always the first item of the lists from $G$.

For the Unique Website Model, we have

$$S_{e|c} = \sum_{s \in Sites(c)} \frac{1}{\sqrt{AvgRank_{c,e,s}}}$$

### B. Collaboative Filtering

For collaborative clustering, the Pearson correlation based similarity and adjusted cosine similarity methods are used.
Using the linear combination of these methods, user can get predicted results.

1) Pearson correlation-based Similarity

$$sim(k, l) = \frac{cov(k, l)}{\sigma_k \sigma_l}$$
$$= \frac{\sum_n^{u,k} (R_{u,k} - \bar{R}_k)(R_{u,l} - \bar{R}_l)}{\sqrt{\sum_{u=l}^n (R_{u,k} - \bar{R}_k)^2} \sqrt{\sum_{i=l}^n (R_{u,l} - \bar{R}_l)^2}}$$

Where, $sim(k, l)$ means the similarity between item $k$ and $l$, $n$ means the total number of users, who rated on both item $k$ and $l$, $\bar{R}_k, \bar{R}_l$ are the average ratings of item $k$ and $l$, respectively; $R_{u,k}$, $R_{u,l}$ mean the rating of user $u$ on item $k$ and $l$ respectively.

2) Adjusted Cosine Similarity

$$sim(k, l) = \frac{\sum_{u=l}^n (R_{u,k} - \bar{R}_u)(R_{u,l} - \bar{R}_u)}{\sqrt{\sum_{u=l}^n (R_{u,k} - \bar{R}_u)(R_{u,l} - \bar{R}_u)^2} \sqrt{\sum_{i=l}^n (R_{u,k} - \bar{R}_u)(R_{u,l} - \bar{R}_u)^2}}$$

Linear Combination $sim(k, l) = sim(k, l)_{item} \times (1 - c) + sim(k, l)_{group} \times c$

Where, $c$ Means the combination coefficient, $sim(k, l)_{item}$ Means that the similarity between item $k$ and $l$,
$sim(k, l)_{group}$ Means that the similarity between item $k$ and $l$

Collaborative Prediction Prediction for an item is then computed by

$$P_{u,k} = \bar{R}_k + \frac{\sum_{i=l}^n (R_{u,i} - \bar{R}_i) \times sim(k, i)}{\sum_{i=l}^n |sim(k, i)|}$$

## V. ANALYSIS AND RESULTS

### A. Dataset

For evaluation, product data is collected from web having different categories such as gender, brand, colors etc.

### B. Expected results

To evaluate the performance, time require to search query on various database size is used. The expected results are evaluated according to time requires to process user query1, query2 and query 3 that extract output result. The table 1 show readings for processing time require for query processing.

**Table.1. Time Require For Search Result**

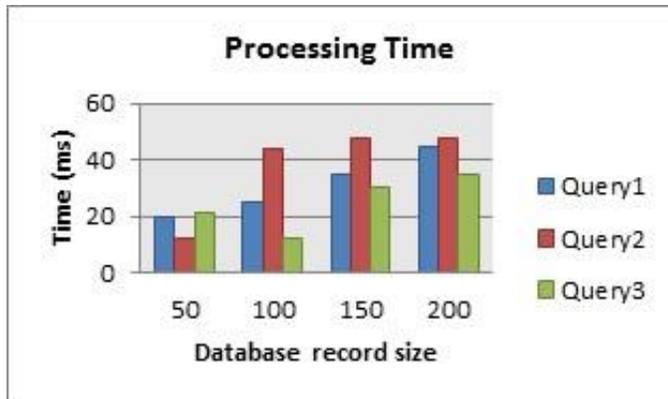| Database size | Query1 | Query2 | Query3 |
|---|---|---|---|
| 50 | 20 | 12 | 21 |
| 100 | 25 | 44 | 12 |
| 150 | 35 | 48 | 30 |
| 200 | 45 | 48 | 35 |



**Figure. 2. Time require for search result**

In item based collaborative method, which makes prediction only based on item-based matrix as in table 2, it is impossible to make predictions on this item. ICHM matrix presentation of some of our products and users are shown in figure 3.
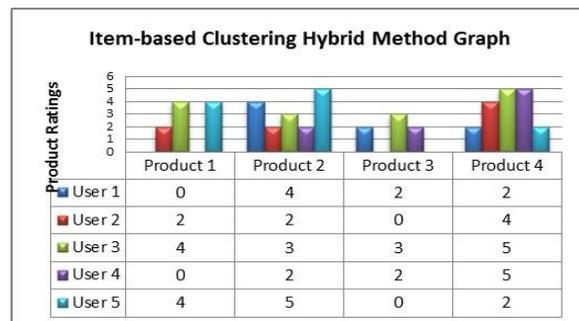


**Figure. 3. ICHM Matrix Presentation**

In user based collaborative method, which can makes prediction for users, based on group rating. UCHM matrix presentation of some of our products and users are shown in figure 4.
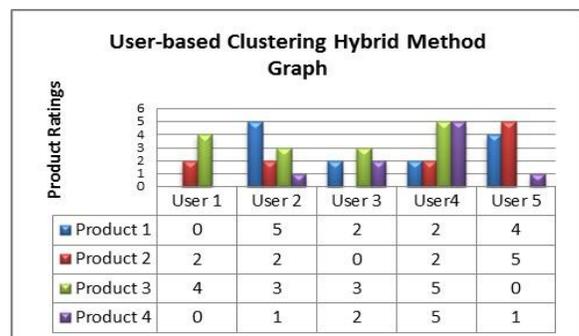


**Figure. 4. UCHM Matrix Presentation**

## VI. CONCLUSION

This paper suggests the query facet which is a collection of items which describe and summarize significant aspect of a query. This paper address the issue of finding query facets which are numerous groups of words or else phrases that clarify and summarize the content enclosed by a query. Paper assume that the significant aspects of a query are generally accessible and frequent in the query's top-k retrieved documents in the form of lists, as well as query facets can be extracted by aggregating these important lists. A systematic resolution is suggested which denoted as QDMiner, to automatically extract query facets by extracting as well as grouping repeated lists from HTML tags, free text and duplicate regions within topk search results. Moreover this paper introduces clustering techniques to the item content data to accompaniment the user rating statistics, which increases the accuracy of collaborative similarity. Using collaborative idea, effectiveness of scheme get increase since of user intends documents is recommended to user. Therefore user search time becomes less for the same or similar data that require to user

## VII. ACKNOWLEDGMENT

## VIII.REFERENCES

[1].L. Bing, W. Lam, T.-L. Wong, and S. Jameel, Web query reformulation via joint modeling of latent topic dependency and term context, ACM Trans. Inf. Syst., vol. 33, no. 2, pp. 6:16:38, eb. 2015.

[2].W. Kong and J. Allan, Extending faceted searchto the general web, in Proc.ACMInt. Conf. Inf. Knowl. Manage., 2014, pp. 839848.

[3].I. Szpektor, A. Gionis, and Y. Maarek, Improving recommendation for long-tail queries via templates, in Proc. 20th Int. Conf. World Wide Web, 2011, pp. 4756.

[4].X. Xue and W. B. Croft, Modeling reformulation using query distributions, ACM Trans. Inf. Syst., vol. 31, no. 2, pp. 6:16:34, May 2013.

[5].L. Li, L. Zhong, Z. Yang, and M. Kitsuregawa, Qubic: An adaptive approach to query-based recommendation, J. Intell. Inf. Syst., vol. 40, no. 3, pp. 555587, Jun. 2013.

[6].W. Kong and J. Allan, Extracting query facets from search results, in Proc. 36th Int. ACM SIGIR Conf. Res. Develop. Inf. Retrieval, 2013, pp. 93102.

[7].Qing Li and Byeong Man Kim, An Approach for Combining Contentbased and Collaborative Filters, Korea Research Foundation Grant (KRF2002-041-D00459), 2002.

[8].I. Szpektor, A. Gionis, and Y. Maarek, Improving recommendation for long-tail queries via templates, in Proc. 20th Int. Conf. World Wide Web, 2011, pp. 4756.

[9].J. Pound, S. Paparizos, and P. Tsaparas, Facet discovery for structured web search: A query-log mining approach, in Proc. ACM SIGMOD Int. Conf. Manage. Data, 2011, pp. 169180.

[10].M. Diao, S. Mukherjea, N. Rajput, and K. Srivastava,, Faceted search and browsing of audio content on spoken web, in Proc. 19th ACM Int. Conf. Inf. Knowl. Manage., 2010, pp. 10291038.

[11].K. Balog, E. Meij, and M. de Rijke, Entity search: Building bridges between two worlds, in Proc. 3rd Int. Semantic Search Workshop, 2010, pp. 9:19:5.

[12].M. Bron, K. Balog, and M. de Rijke, Ranking related entities: Components and analyses, in Proc. ACM Int. Conf. Inf. Knowl. Manage., 2010, pp. 10791088.

[13].C. Li, N. Yan, S. B. Roy, L. Lisham, and G. Das, Facetedpedia: Dynamic generation of query-dependent faceted interfaces for wikipedia, in Proc 19thInt.Conf World Wide Web, 2010, pp.651660 2010, pp. 651660.

[14].A. Herdagdelen, M. Ciaramita, D. Mahler, M. Holmqvist, K. Hall, S. Riezler, and E. Alfonseca, Generalized syntactic and semantic models of query reformulation, in Proc. 33rd Int. ACM SIGIR Conf. Res. Develop. Inf. retrieval, 2010, pp. 283290.

[15].J. Huang and E. N. Efthimiadis, Analyzing and evaluating query reformulation strategies in web search logs, in Proc. 18th ACM Conf. Inf. Knowl. Manage., 2009, pp. 7786.

[16].S. Gholamrezazadeh, M. A. Salehi, and B. Gholamzadeh, A comprehensive survey on text summarization systems, in Proc. 2nd Int. Conf. Comput. Sci. Appli., 2009, pp. 16.

[17].H. Zhang, M. Zhu, S. Shi, and J.-R. Wen, Employing topic models for pattern-based semantic class discovery, in Proc. Joint Conf. 47th Annu. Meet. ACL 4th Int. Joint Conf. Natural Lang. Process. AFNLP, 2009, pp. 459467.

[18].O. Ben-Yitzhak, N. Golbandi, N. HarEl, R. Lempel, A. Neumann,S. Ofek-Koifman, D. Sheinwald, E. Shekita, B. Sznajder, and S. Yogev, Beyond basic faceted search, in Proc. Int. Conf. Web Search Data Mining, 2008, pp. 3344.

[19].W. Dakka and P. G. Ipeirotis, Automatic extraction of useful facet hierarchies from text databases, in Proc. IEEE 24th Int. Conf. Data Eng., 2008, pp. 466475.

[20].S. Riezler, Y. Liu, and A. Vasserman, Translating queries into snippets for improved query expansion, in Proc. 22nd Int. Conf. Comput. Ling., 2008, pp. 737744.