

# International Journal of Innovative Research in Science, Engineering and Technology

(A High Impact Factor, Monthly, Peer Reviewed Journal)

Visit: [www.ijirset.com](http://www.ijirset.com)

Vol. 7, Issue 7, July 2018

## Improving Web Search by Mining Query Facets from Search Results

Gunjan H. Deshmukh<sup>1</sup>, Govinda Borse<sup>2</sup>

M. E Student, Department of Computer Science and Engineering, G.H.Raisoni Institute of Engineering and  
Management, Jalgaon, India<sup>1</sup>

Assistant Professor, Department of Computer Science and Engineering, G.H.Raisoni Institute of Engineering and  
Management, Jalgaon, India<sup>2</sup>

**ABSTRACT:** Query Faceted search is a way for searching to find, analyse, and navigate through search data from online web pages. An effective approach for faceted search is the scope of this implementation. Most existing faceted search and facets generation systems are built on a specific domain (such as product search) or predefined facet categories. For example, Web search mining for an unsupervised contents by automatic extraction of facets that are relevant for search result for personal web search as user search interest pattern from text databases. Facet hierarchies are generated for a whole collection, instead of for a given query. Proposed facets searching system for information discovery and media exploration in online search results. The proposed system explore to automatically find query related aspect of search for open-domain queries in Web search engine. Facets of a query are automatically mined from the top web search results of the query without any additional domain knowledge required. As query facets are good summaries of a query and are potentially useful for users to understand the query and help them explore information, they are possible data source that enable a general open-domain faceted exploratory search. In order to increase the quality of the extracted list, wrapper data extraction concept is introduced.

**KEYWORDS:** Facets generation, Query aspects, Query Reformation, Seed Collection

### I. INTRODUCTION

The plentiful content of the World-Wide Web is useful to millions. In this case, users submit a query, typically a list of keywords, and receive a list of Web pages that may be relevant, typically pages that contain the keywords. There are many challenges in building good search engines, and describe some of the techniques that are useful. Query aspects provide interesting and useful information about a query and, therefore, can be used to improve research experiences in many ways. First, we can show the facets of the query with the innovative results given by search engine appropriately. Therefore, users can understand some important facets of a query without go through number of pages. In this paper, a proposed system scans to automatically identify the look related to searching for open domain queries in the Web search engine. The facets of a query are directly extracted from the consequences of the main query web search lacking the requirement for further domain knowledge. Because the aspects of the consultation are summaries in well manner and are truly helpful for users to know the query and aid them to get the information, data sources are possible that allow a general multifaceted exploratory search of open domain. System address the problem of finding query facets and to link to the intended web page. A query facet is a group of items which outline and encapsulate one prime aspect of a query. The proposed system address the difficulty of finding query facets which are several groups of words or phrases that explain and sum up the content covered by a query. This system bases on open domain search, it is not related only specific people or product. We are using online datasets for getting good results so for any type of query the system will give more accurate and most relevant result to the user. Query dimensions provide interesting and useful knowledge about a query and thus can be used to improve search experiences in many ways. First, we can

# International Journal of Innovative Research in Science, Engineering and Technology

(A High Impact Factor, Monthly, Peer Reviewed Journal)

Visit: [www.ijirset.com](http://www.ijirset.com)

Vol. 7, Issue 7, July 2018

display query dimensions together with the original search results in an appropriate way. Thus, users can understand some important facets of a query without browsing tens of pages.

## II. LITERATURE REVIEW

- In this thesis author consider novel semantic presentation for query subtopic is implemented, which covers phrase embedding approach and query classification distributional representation, to solve those problems mentioned above. Additionally this approach combines multiple semantic presentations in the vector space model and calculates a similarity for clustering query reformulations. Furthermore, automatically discover a collection of subtopics from a query given by user and each of them are presented as a string that define and disambiguates the search intent of the original query. Query subtopic could be mined from various resources involving query suggestion, top-ranked search results and external resource [2].
- In this article, author represents concept of query aspect to know the user interest for search on various topics, where every facet presents a collection of words which explain an user intention for searching that a query. Investigated approach generates subtopics related to query factors and proposed faceted diversification approaches. The actual query aspects are investigated to give more specific search to user such as collecting facets and exploratory search. [1]
- In this thesis, author presents OLAP model for online analysis of user interest to extract query aspects with OLAP capabilities, existence of facet mining was related to the data over relational database, to the free text queries from metadata list style content. This is an extension shows efficiently facet extraction by a search engine to support correlated facets a more difficult data model in which having the values related with searches in numerous facets which are not autonomous [5]
- In this survey author proposes a random faceted search approach for searching query driven analysis on data with both textual content and structured attributes. From a facet query, user expected to dynamically choose a set of interesting aspect and present to a user. Similarly in OLAP exploration, author defines interestingness as how collected value is, based on a given probability [6]. Author of this term paper searched the new techniques based on a graphical representation to extract query facets from the collection of noisy data. The graphical representation tells how likely a candidate form is to be a aspect string and how two terms are clustered jointly in a query feature, and captures the dependency between the two factors. This work proposes two mechanism for aggregation of an inference on the graphical representation since exact deduction is obdurate [4].
- A hidden webpage extraction from an association makes easy to get to the maze by allowing end user to enter queries by a search engine. In other way, data collection as of such a cause is not by implemented in hyper link. As an alternative, data are collected by querying the boundary, and understanding the consequence page with randomly generated [3].
- This paper resolve problem of relevant search by using the stuffing of pages to focus the search on a subject; by prioritizing capable links within the topic; and by also following links that may not lead to instant advantage. This paper recommend a new techniques whereby searching involuntarily find out patterns of useful links and apply their spotlight as the creep progresses, thus mainly plummeting the amount of requisite physical setup and change [8].
- In this paper author design a two-stage crawler, namely Smart Crawler, for related harvesting deep web pages. In the primary stage, Smart Crawler performs web site (URL) based penetrating for hidden web pages with the assistance of search engines, avoiding visiting a huge quantity of pages. To accomplish extra efficient outcome for a focused crawl, Smart Crawler position webpage to prioritize extremely related data for a specified search query. In the next stage, Smart Crawler achieves rapid in site web crawling by extracting mainly related links with an adaptive link prioritizing [7].

# International Journal of Innovative Research in Science, Engineering and Technology

(A High Impact Factor, Monthly, Peer Reviewed Journal)

Visit: [www.ijirset.com](http://www.ijirset.com)

Vol. 7, Issue 7, July 2018

## III. SYSTEM ARCHITECTURE

Query facets provide fascinating and useful knowledge about a query and thus can be used to upgrade search experiences. Users can perceive some prime aspects of a query without browsing tens of pages. Query facets may provide direct information or instant answers that users are seeking. The facets are generated based on the user's taste. The user can also go to the desired web page by selecting the item from the facets. Thus it saves the user's time wasted on searching for the data in tens or thousands of pages. Query facets are mined by using four methods. Processing steps are as follows:

### 1. Crawling (Seed Collection):

A web crawler forms an integral part of any search engine. The basic task of a crawler is to fetch pages, parse them to get some more URLs, and then fetch these URLs to get even more URLs. In this process, a crawler can also log these pages or perform several other operations on pages fetched according to the requirements of the search engines. The crawler module retrieves pages from the Web for later analysis by the indexing module.

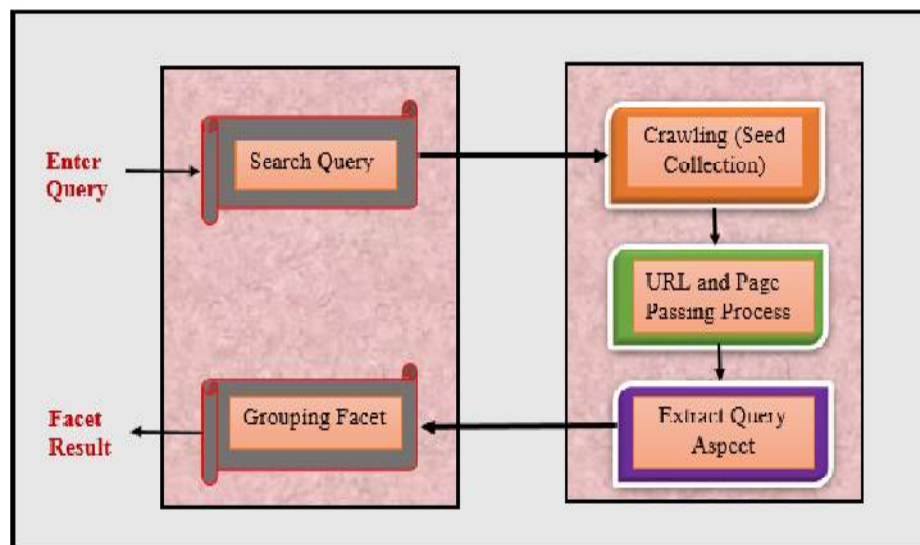


Fig.1 System architecture of facet mining

### 2. URL and Page parsing process:

First system consider URL, it will take last word from URL as a facet also for a list extracted from a HTML element like SELECT, UL, OL, or TABLE by pattern. That contain facet and links that will display to user.

### 3. Extract Query Aspects:

After performing page extraction we get facets and links. SELECT For the SELECT tag, we simply extract all text from their child tags (OPTION) to create a list. UL/OL For these two tags, we also simply extract text within their child tags (LI). For a list extracted from a HTML element like SELECT, UL, OL, or TABLE by pattern HTMLTAG, its context is comprised of the current element and the previous and next element if any. Facets mainly extracted from the URL link and title.

# International Journal of Innovative Research in Science, Engineering and Technology

(A High Impact Factor, Monthly, Peer Reviewed Journal)

Visit: [www.ijirset.com](http://www.ijirset.com)

Vol. 7, Issue 7, July 2018

## 4. Facet Grouping:

Facets are grouped according to different classes. It group data of similar facetssould frequently appear in the top results. This model isapplied over html document by parsing html tags. List clustering Similar lists are grouped together to compose a facet.

## IV. METHODOLOGY

Query facets can be used to help improve search experience in many ways. Like in facetedsearch [6], query facet can help users to navigate through different topics of the search resultsby applying multiple filters.

Then, query facets are mined by the following steps:

### 1. Search Query :

As if user want to search related to anything user has to put the Query or keyword respectively,after that user has two options to search regarding that query,either user can search itin generalize way or in personalize way.

1) *Generalize Search*:In the generalize search system will take the query from user and send it to the search engine because it uses online database, so it will get the result or this process will be perform as per the normal search.

2) *Personalize Search*:In personalize search,system will give result by considering both query and user interest.Here, at the time of user registration user should put the profession or area of interest.This information later use for personalize search. For example if user put the queryrelated to shoes,and user area of interest is Sports so in personalize search user considerboth shoes and sports and in the result it will return the links related to correspondingquery and user interest,So user get most of links and facets related to the sports shoesbrands.

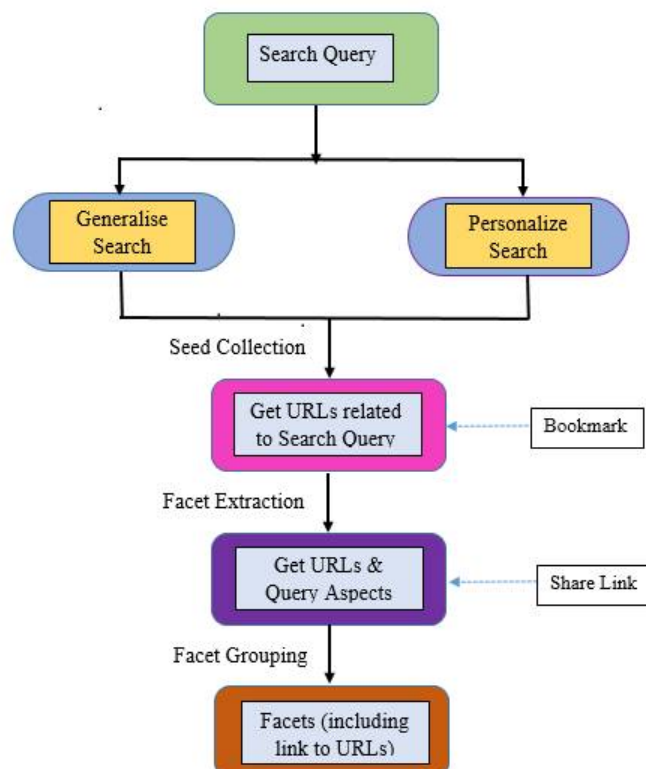


Fig.2 working of proposed system

# International Journal of Innovative Research in Science, Engineering and Technology

(A High Impact Factor, Monthly, Peer Reviewed Journal)

Visit: [www.ijirset.com](http://www.ijirset.com)

Vol. 7, Issue 7, July 2018

## 2 . URLs related to Search Query:

After performing first step of entering and submitting search query seed collection operation is performed. In the seed collection step search engine will search all the webpages or links related to the search query. Firstly, it will take the raw query and remove the stop words from it and remaining keyword is sent to the search engine. System will return the webpages and URLs related to user's query in the form of list to the user. In this step system provides the option of Bookmark to user, in this step user can bookmark the important link for further use. Bookmark links are easily available to user for quick access so it will reduce the searching time and efforts.

## 3. Query Aspects and URLs:

In the next step, facet extraction process is performed in this step the system will find out different dimensions related to the search query. The query dimensions are extracted from webpages and different URL links or URL title and shown to the user in the form of list. In this step system provides option of Share Link with the help of this option user can send the link to the friends. Hence, due to this process unnecessary recommendation will be avoided. If another user wants to search related to the same query so with share link option important links can be shared to that user, so instead of visiting number of links only main link will be visited.

## 4. Query Facets:

In this step, before displaying the facets to user Facet grouping step will be performed. In facet grouping step the collection of query aspect is performed and after grouping of facet it will show to the user. The group of facets includes all the facets related to the search query, these facets contain link to the URLs related to that facets. After clicking on that facets it will show links of different webpages. In this along with the main facets user can also parse or visit another extra facets.

## V.EXPERIMENTAL RESULT

Proposed system mine and analyse query facets or query aspects from search query. After the final result is produced, all facets are displayed to the user and at that time classification graphs are generated.

### 1. Generalize search :

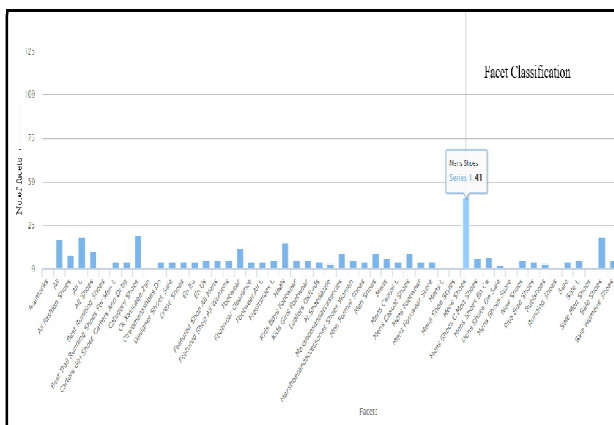


Fig 3. a) Facet Classification (Generalize)

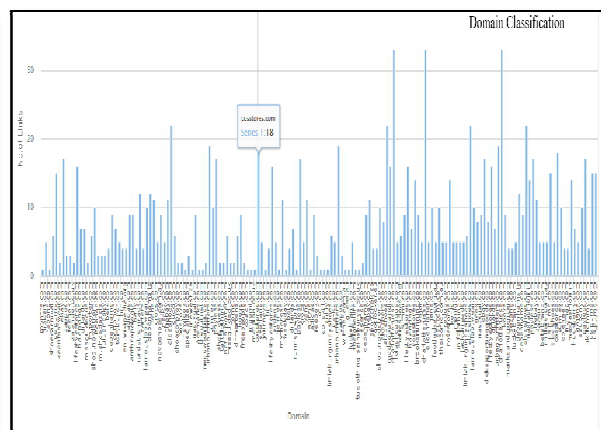


Fig 3. b) Domain Classification (Generalize)

# International Journal of Innovative Research in Science, Engineering and Technology

(A High Impact Factor, Monthly, Peer Reviewed Journal)

Visit: [www.ijirset.com](http://www.ijirset.com)

Vol. 7, Issue 7, July 2018

Fig 3 a) shows facet classification analysis graph, In this the user put query for searching and search it in generalize way for example user put query 'Shoes' and search it in normal way it will give all aspects related to the query shoes and in this process only that query is consider as input not user interest nor profession is include. In the output the group on facets along with URL links are display to user. b) This graph represents, how many links contains by each domain, so user can see domain wise classification of facets. If particular domain includes more links then user can visit that domain or the URLs contain by that domain first. Domain classification graph shows the proper distribution of domain and the links under that particular domain.

## 2. Personalize Search:

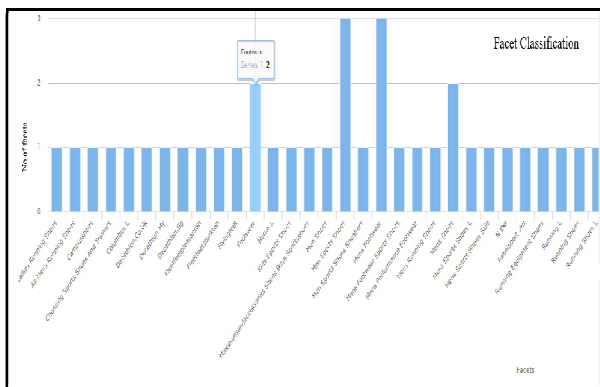


Fig 4. c) Facet Classification (Personalize)

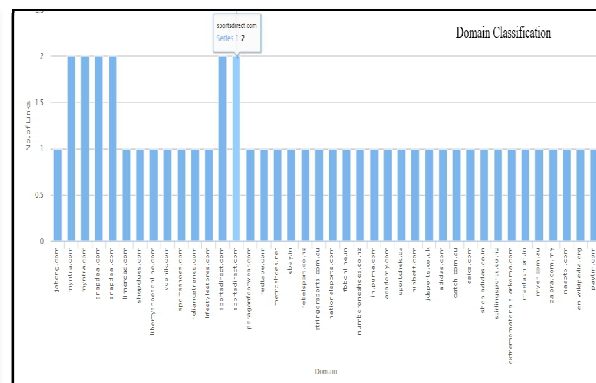


Fig 4. d) Domain Classification (Personalize)

Fig 4 c) shows the analysis graph of facets of search query. In facet classification in personalize case first when user put query for the search then user should select the way to search that particular query. As if user select the personalize way then while sending the search query to search engine, user interest also get consider. This user interest is consider from the user profile, because at the time of registration user put the all information along with area of interest. If user choose the personalize way to search particular query then user's area of interest and search query both gets consider. For example at the time of registration user puts area of interest as 'Sports', and at the time of searching user put query as 'Shoes' and select the option of personalize search then the system will return more links related to the sports shoes. So it will reduce the searching efforts and within less time user get the proper result for search query. b) Represents domain classification graphs. When system gives number of links as a output, this analysis shows the graphical representation of domain. As if search for user's query is perform in personalize way it will consider or shows only that domain where system finds aspects by considering user query and user's area of interest both.

## VI. APPLICATION

In web application development, educational website, e-learning applications, online literature needs to generate facet for better search results. It is useful for getting the results of queries in less time. It will save the time by providing the answer for the given queries along with some extra facets. If user wants to see another aspect of the queries, so without hitting another sites, this system will give the information for that aspect in the same result. Hence the unnecessary recommendation is avoided for other user. This concept of facets mining and analysis will be applicable in data mining. Users get most relevant result related to their search queries in less time, many recommendation system application can use it.



# International Journal of Innovative Research in Science, Engineering and Technology

(A High Impact Factor, Monthly, Peer Reviewed Journal)

Visit: [www.ijirset.com](http://www.ijirset.com)

Vol. 7, Issue 7, July 2018

## VII. CONCLUSION AND FUTURE SCOPE

Hence we studied and implemented the method of generating meaningful facets from the user query search results. The facets here are generated using four steps, web crawling, URL parsing, facet extraction and facet grouping. In web crawling, user send query to search engine and it will find relevant result related to that query. In URL Parsing the unique websites are collected and shown in the list. This method can extract high quality lists from the top k query search results. Hence these high quality lists can be used to generate meaningful facets. In facet extraction step the facets are extracted containing links. After this process grouping of facets is performed. These facets are generated based on either user's interest that means in a personalized way or in a generalized way. User can navigate to the specified page by selecting the item on the facet to get detailed information. Mining information in the form of facets from record by parsing HTML labels from the report.

In future will also investigate the concept that, along with facets the good descriptions of query facets may be helpful for users to better understand the facets. Automatically generate meaningful descriptions is an interesting research topic, this will be developed for e-learning application. Faceted search is a technique for accessing information. Mining facets classification system allows users to collect information by multiple filters. Finding query facets is an important issue due to this only text based query processing is used. If image is used as a query this makes complicated one so in future will investigate the concept for better result.

## REFERENCES

- [1] M. Damova and I. Koychev, "Query-based summarization: A survey," in Proc. S3T, 2010, pp. 142–146.
- [2] K. Shinzato and T. Kentaro, "A simple www-based method for semantic word class acquisition," in Recent Advances in Natural Language Processing (RANLP '05), pp. 207–216, 2005.
- [3] 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. 459–467.
- [4] Yassin M. Y. Hasan and Lina J. Karam, "Morphological Text Extraction from Images", IEEE Transactions On Image Processing, vol. 9, No. 11, 2000
- [5] M. J. Cafarella, A. Halevy, D. Z. Wang, E. Wu, and Y. Zhang, "Webtables: Exploring the power of tables on the web," VLDB, vol. 1, pp. 538–549, Aug. 2008.
- [6] J. K. Latha, K. R. Veni, and R. Rajaram, "Afgf: An automatic facet generation framework for document retrieval," in Proc. Int. Conf. Adv. Comput. Eng., 2010, pp. 110–114
- [7] Aria Pezeshk and Richard L. Tutwiler, "Automatic Feature Extraction and Text Recognition from Scanned Topographic Maps", IEEE Transactions on geosciences and remote sensing, VOL. 49, NO. 12, 2011
- [8] J. S. Basu Roy, H. Wang, G. Das, U. Nambiar, and M. Mohania, "Minimum-effort driven dynamic faceted search in structured databases," in Proc. ACM Int. Conf. Inf. Knowl. Manage., 2008, pp. 13–22.
- [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. 169–180.
- [10] E. Stoica and M. A. Hearst, "Automating creation of hierarchical faceted metadata structures," in Proc. Human Lang. Technol. Conf., 2007, pp. 244–251.
- [11] J. O. Etzioni, M. Cafarella, D. Downey, S. Kok, A.-M. Popescu, T. Shaked, S. Soderland, D. S. Weld, and A. Yates, "Web-scale information extraction in knowitall: (preliminary results)," in Proc. 13th Int. Conf. World Wide Web, 2004, pp. 100–110.
- [12] 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. 466–475.
- [13] D. Dash, J. Rao, N. Megiddo, A. Ailamaki, and G. Lohman, "Dynamic faceted search for discovery-driven analysis," in ACM Int. Conf. Inf. Knowl. Manage., pp. 3–12, 2008.
- [14] W. Dakka and P. G. Ipeirotis, "Automatic extraction of useful Facet hierarchies from text databases" in Proc. IEEE 24th Int. Conf. Data Eng., 2015, pp. 466–475.
- [15] M. Diao, S. Mukherjee, 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. 1029–1038.