

# AUTOMATIC FACET MINING AND RECOMMENDATION FROM WEB SEARCH IN CLOUD ENVIRONMENT

**Kirti Parmar, Minal L. Solanki, Priya Bhat, Sayali Ghorpade**

**Prof. Sneha Farkade**

*Department of Computer Engineering, G.S.Moze College of Engineering, Balewadi, Pune*

## ABSTRACT

*In search engines, different users may seek for different information by issuing the similar query. To convince more users with partial search results, search result diversification re-ranks the results to coat as many user intents as probable. Most presented intent-aware diversification algorithms differentiate user intentions as subtopics, every of which is typically a word, a phrase, or a piece of clarification. Web search queries are often uncertain or multi-faceted, which makes an effortless ranked list of results insufficient. To help information finding for such faceted queries, system explore a technique that explicitly represents fascinating facets of a query using groups of semantically related terms retrieved from search results. As an example, for the query baggage allowance, these groups might be many airlines, different flight category (domestic, international), or multiple travel classes (first, business, economy). System identifies such groups query facets and the conditions in these groups facet classes. In the proposed work system denote a supervised approach based on a graphical model to identification query facets from the noisy candidates found. The identical or visualize model learns how likely a contestant 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. System proposes two algorithms for estimated inference on the visualization model since exact inference is intractable. System's estimate combines accuracy of the facet conditions with the grouping quality. Investigational results on a illustration of web queries show that the supervised approach significantly outperforms existing methods, which are mostly unsupervised, recommending that query facet retrieval can be effectively learned. The evaluation has done with public cloud environment.*

**Keywords:** *Faceted Web Search, Interactive Feedback, Query Facets, Summarization, User Intent*

## I. INTRODUCTION

For most people, the way they interact with web search engines has not changed significantly in the last decade. They still issue queries manually and review lists of result documents. The most significant and obvious user interface changes were the introduction of verticals (e.g. images, videos, and news), query auto complete, and question answering (e.g. Google Knowledge Graph). However, most internet users are also acquainted with faceted search: any ecommerce website, any library and most catalogues of any kind employ this technique to

provide an accessible and fast way to locate arbitrary objects. System believes that most users would appreciate the utilization of this idea in web search. However, this is no trivial task. The ultimate goal of Faceted Web Search is to support the user to accomplish his search task. Previous work focused on the idea of using existing taxonomies or on generating facets for an entire corpus offline after indexation. These approaches lack the adaptation to the document result space or the user intent, and are too narrow. System propose web search facets that automatically recognize different subtopics, partition the search result space evenly and exhaustively per subtopic, and still contain only a small number of terms.

However, these original facets cannot be directly adopted as subtopics. Since query facets are designed for splitting different facets of a query, they are usually far more fine-grained than traditional subtopics in diversification. There are some objectives to achieve the successful implementation of system.

1. Need to impalement novel base recommendation approach suing user feedback or search sessions.
2. To extract the data from the search engine databases related to query, searched by the user and represent the search results in restructured manner.
3. To provide search results according to search goals of particular user.
4. System has to provide a service recommendation base on similarity score which is calculating using text similarity algorithm.
5. Improve the system accuracy using the clustering algorithm base on potential users.
6. Using hierarchical association algorithm finds the user interest and collaborate the filter clusters.

## **II. RESEARCH BACKGROUND**

A great deal of examination works concentrates on database interfaces which help clients to question the social database without SQL Query-By-Example and Query Form are two most generally utilized database questioning interfaces. Current considers and works predominantly concentrate on the most proficient method to create the question shapes.

### **2.1 Adjusted Query Form**

Your devices gives by the database customers make awesome endeavors to offer engineers some assistance with generating the question frames, for example, Easy Query, Cold Fusion etc. They give visual interface to engineers to create or squeeze question frames. The issue of those instruments is that, they are accommodated the specialist engineers. H.V. Jagadish anticipated a system which grants end-customers to change the present inquiry structure at run time. If the database outline is immeasurable, it is troublesome for end customer to discover reasonable database substances and attributes.

### **2.2 Computerized Creation of Forms**

M. Jayapandian exhibited information driven technique. It to begin with finds an arrangement of information properties, which are doubtlessly questioned in light of the database pattern and information examples. At that point, the inquiry structures are created taking into account the chose characteristics.

### **2.3 Mechanizing the configuration and development of question structures**

H.V. Jagadish introduced a workload-driven strategy. It applies grouping calculation on recorded questions to discover the agent inquiries. The question structures are then created in light of those agent questions. One issue of the previously stated methodologies is that, in the event that we create bunches of inquiry structures ahead of

time, there are still client inquiries that can't be fulfilled by any of question structures. Another issue is that, when we produce an extensive number of question shapes, how to let clients locate a suitable question structure would be testing.

#### 2.4 Joining catchphrase look and shapes

An answer for earlier stated methodologies is projected in. It accordingly produces a substantial measure of question shapes ahead of time. The client inputs a few catchphrases to discover significant question frames from a significant number of pre-created difficulty shapes yet it is not suitable when the client does not have solid watchwords to depict the investigation

### III. LITERATURE SURVEY

Comprising single or more facet conditions in order to verify, whether the compound term selected by a user is suitable. The second function select only a section of all the facets to display in the user interface, when there are too numerous applicable facets and facet terms. The query refinement process runs as follow. First, the facets in a faceted classification are ranked, and some are chosen to be showed in the user interface. The terms in the selected facet are working by user to form composite terms. Second, the system validates the compound terms selected by the user, and show only the data objects associated with the valid combined terms. At the same time, the system updates the number of data objects subsequent to facet terms in user interface and ranks the facet terms again for the next navigation activity. The iterations continue until expected outcome are establish. Search outcome ranking in the faceted search is like to that in the conventional Information recovery domain. It has been extensively studied for years [3-5]. Thus we will skip it from the following sections. Manual recognition of facet terms by domain experts is costly, wasteful and has poor scalability. There is a few researchers have conducted a preface study of automatic facet term pulling out. Existing automatic removal methods can be divided into three categories corresponding to the different data types: unstructured, semi-structured and structured.

Unstructured data refers to the information that does not adhere to a predefined data model. In facets term extraction, the most regular form of formless data is the natural language text, which is always uncertain and ill-formed. Since machine understanding of normal language text ruins as an open subject, it is very hard to automatically extract facet terms by machine learn techniques only. Current removal methods mainly focus on making comprehensive use of the information of terms, linguistic features of the terms and external knowledge base. The typical methods are outlined as follows.

**Stoica et al. [6]** proposed Castanet algorithm to select facet terms based on term frequency distribution. The essence of this algorithm is selecting the conditions having an occurrence higher than a threshold as facet term candidates for following processing. This algorithm can be easily implemented and extended to different domains since only term frequency is employed. **Anick and Tipirneni [7]** proposed a facet term extraction algorithm based on the lexical dispersion of words in text. Lexical distribution of a word is the amount of different compounds that contain this word within a document group. The algorithm consists of two stages. In the indexing stage documents are parsed so the lexical compounds can be extracting. In the querying stage the compound appear in the top  $n$  documents of a ranked result list are used to compute the lexical distribution of each term happening within these compounds. These terms are then sorted by their dispersions and also the top

$m$  conditions are chosen as candidate terms for following hierarchy construction. The disadvantage of this algorithm is that the removal of facet conditions depends on the specific lexical structure and so can be hardly extended to new domains.

**Ling et al. [8]** proposed a two-stage probabilistic method to extract facet terms based on topic model. Given the original keywords from a user, these techniques first apply a bootstrapping algorithm to the document collection to get more correlated terms. Probabilistic combination models are applied to these extended terms to estimate the term distribution of every facet. This is done by at the same time fit the topic model to the data set and restraining the model so that it is secure to the specified definition from the user. The basic thought behind the processes is to guide the subject model with user-defined keywords.

Dakka and Ipeirotis proposed an unsubstantiated automatic facet removal algorithm using exterior resources [9]. This algorithm first identifies the facet term candidates in each document by using third-party term extraction services or algorithms, such as **Ling Pipe [10]**, **Yahoo Term Extraction [11]**, Wikipedia, or **Taxonomy Warehouse [12]**. Then, each candidate is extended with "context" phrases appearing in external assets by querying WorldNet, Wikipedia, and other online dictionary. This step produces the latent facet terms in the expanded term set, which do not explicitly come out in the documents. Finally, the term distributions in the new term set and the extended term set were compared to identify the terms that can be used to assemble browsing facets. This algorithm has good flexibility and extensibility. However the value of the extracted facets greatly depends on the value of the external resources and term extractor.

The partially structured data does not conform to an unambiguous data diagram; however, it generally contains tags or other markers to divide semantically related essentials. Examples of the semi structured data include HTML pages and the pages annotate by Resource Description Framework (RDF). Partly structured data has an implicit formal structure, which can be exploited to recover the quality of facet term removal. For example, the hyperlinks of web pages can be used to evaluate the magnitude of facet terms. The typical removal methods in this subfield are described briefly as follows.

**Roy et al. [13]** presented such a method. At every step, a user is asked one or more questions about the different facet conditions, and the mainly promising set of facet conditions is identified based on the user's response.

**Zhao et al. [14]** implement Explorer which selects facet terms from the attributes by measuring the significance between documents and keywords.

**Dou et al. [15]** developed QDMiner system to automatically extract facet hierarchies for keywords by aggregating frequent lists from HTML tags, free text, and repeated regions within top search results. The process is described as:

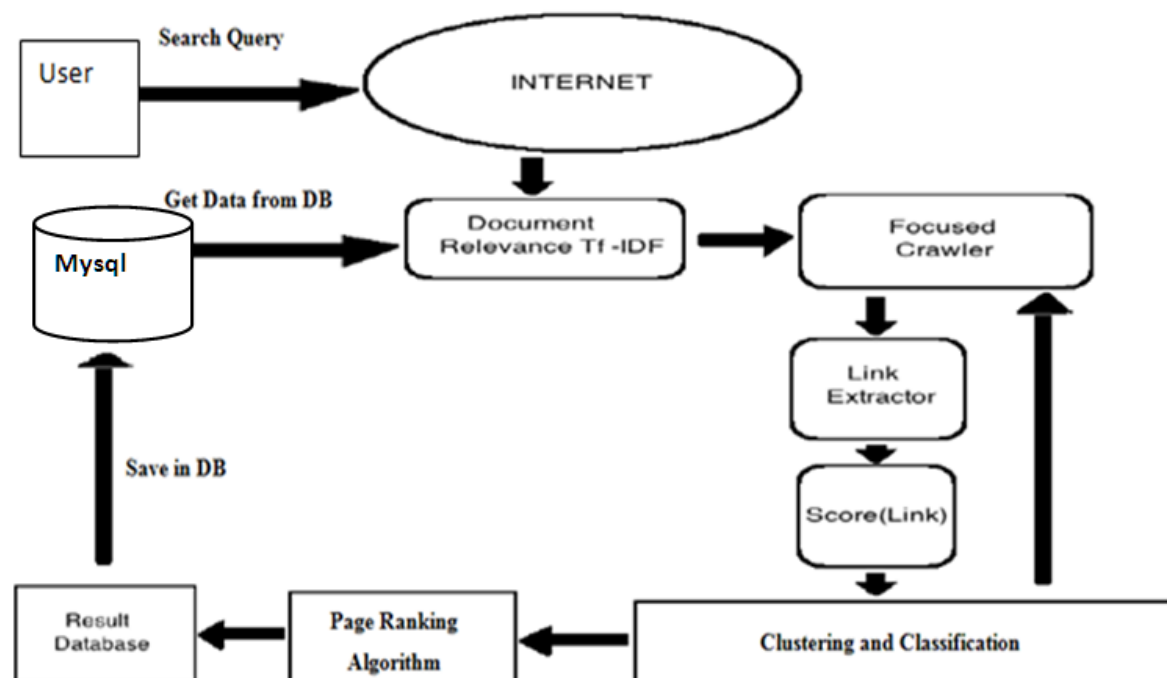
1. When a user issue a query, QDMiner retrieve the top-k outcome from a search engine to form a set  $R$ .
2. QDMiner extracts and weights several types of list from every document of  $R$ .
3. QDMiner groups similar lists together to form a facet by a modified value threshold cluster algorithm.
4. QDMiner evaluates and ranks facets and facet items based on their importance.

**Chen and Li [16]** presented a method of automatic term hierarchies acquisition based on subsumption estimation and spectral clustering. First, every term is considered as a vertex in an undirected weighted graph. The problem of hierarchical relation creation is then modeled as a customized graph-partitioning problem and is

solve by spectral clustering methods. Subsumption estimation is introduced to direct the spectral cluster route. As a result, a modified spectral graph partitioning algorithm was developed to correctly depict the hyponym information about facet terms. This method can extract facet conditions based on composite words such as “probabilistic clustering,” but the hierarchical taxonomies may be significantly different from the physically obtained taxonomies.

**Oren et al. [17]** proposed a facet term recognition method dedicated to the semi-structured data in semantic web. This technique was implementing by dynamically construct a faceted routing tree based on RDF graph. Structured data has an explicit data representation or scheme, such as data store in a relational database. For structured data, the core task of facet term removal is to choose facet terms from attribute of database. He exploited the predicate balance, predicate frequency and object cardinality to rank facets. The predicate balance is referred to as the balance of the faceted direction-finding tree composed of faceted attribute. The more balanced this navigation tree is, the higher the navigation competence is. The object cardinality is the amount of values that can be assigning to the faceted attribute. The lesser this cardinality is, the easier it is for user to select a suitable value. The predicate frequency indicates the categorization capacity of a facet. If the predicate occurrence of a faceted attribute is low, when a user selects a value of this feature, only small quantities of data items are affected.

#### IV. PROPOSED SYSTEM



**Fig. 1: Proposed System Architecture**

In the proposed research work to design and implement to automatically mine query facets by extracting and grouping frequent lists from free text, HTML tags, and repeat regions within top search results with higher

security in cloud framework. Once user submit any query system first check the availability in existing sessions with clustering list. If present query is available in existing sessions database it will return all URL's from database. System will deploy the system on Amazon EC2 console as public cloud. It also focus on database security approach using SQL injection and prevention techniques. Using MVC architecture system provide system boosting for generate minimum load on database as well server side, and improve the system performance. System has deployment on EC2 base public cloud environment. Provide database security using Aho-corasick base SQL injection and prevention algorithms. It store session history of each user into local database it will help for clustering and ranking approach.

#### 4.1 Algorithms

System performs a novel diversity-aware service ranking algorithm to find the optimal top-k Web services based on a proposed comprehensive ranking measure. It is re-lasted work on service recommendation in these three categorists, and on diversity-based ranking algorithms.

##### 4.1.1 Document retrieval Algorithm

**Input:** Users query as Q, Network Connection N;

**Output:** result from relevancy calculation top k pages base on Q.

Step 1: User provide the Q to system.

Step 2: if (N!=Null)

    Process

    Read each attribute A from ith Row in D

    Res[i]=Calcsim(Q,A)

Else No connection

Step 3: For each(k to Res)

Step 4: ArrayList Objarray to bind Q to Res[i] or k

Step 5: Return to users Objarray

Step 6: Display Objarray

##### 4.1.2 Weight Calculation Algorithm

**Input:** Query generated from user Q, each retrieved list L from webpage.

**Output:** Each list with weight.

Here system has to find similarity of two lists:  $\vec{a} = (a_1, a_2, a_3, \dots)$  and  $\vec{b} = (b_1, b_2, b_3, \dots)$ , where  $a_n$  and  $b_n$  are the components of the vector (features of the document, or values for each word of the comment) and the  $n$  is the dimension of the vectors:

Step 1: Read each row R from Data List L

Step 2: for each (Column c from R)

Step 3: Apply formula (1) on c and Q

Step 4: Score=Calc(c, Q)

Step 5: calculate relevancy score for attribute list.

Step 6: assign each Row to current weight

Step 7: Categorize all instances

Step 8: end for end procedure

#### 4.1.3 Mathematical Model

The system has classified into the different sets like below

Sys= {inp, process, out, analysis}

Inp= {Q1, Q2.....Qn}

That is the set of input queries

List= {L1, L2, L3.....Ln}

$$LD[w]: \sum_{k=0}^n \binom{n}{k[D]}$$

Extracted list from each documents

L: {Wi1, Wi2, Wi3..... Win} weight of each list using below formula

$$\vec{a} \cdot \vec{b}: \sum_{i=1}^n a_i b_i = a_1 b_1 + a_2 b_2 + \dots + a_n b_n$$

C= {c1, c2.....cn} clusters of each list

$$[C1...Ck]: \sum_{k=0}^n k(\text{classify})$$

The finally system work with item or facet rank it can create the set of higher dimensions.

UrlList= {URL1 (w), URL2 (w)..... URLn (w)}

$$UL[k]: \sum_{n=1}^m Doc1 + Doc2 \dots \dots \dots Docm$$

Success condition

if(inp! = Null)

Failure condition

if(UrlList == Null)

## V. IMPLEMENTATION DETAILS

### 5.1 System Requirements

1. System interfaces: Windows Operating System
2. User interfaces: User interface using Jsp and Servlet
3. Hardware interfaces

Processor: - Intel Pentium 4 or above

Memory: - 512 MB or above

Other peripheral: - Printer

Hard Disk: - 40GB

### 4. Software interfaces:

Front End: Jdk 1.6.0, Netbeans 6.9.1, IE 6.0/above

Back-End: Mysql 5.1



**5. Communications interfaces**

We will use TCP/IP protocol for establishing connection and transmitting data over the network. We will use Ethernet for LAN

**6. Memory constraints**

Basically initial software's will use around 20 GB on hard drive, and our actual application will be around 300 MB data. When we deploy the application on web server we will assume the 1 GB space for website and 500 MB for database.

**7. Site adaptation requirements**

A domain name, hosting plan on web server and Database space required. Some possibility of using hadoop server base on load balancing

**8. Services: Amazon EC2 for cloud****VI. RESULTS AND DISCUSSIONS**

Comparison between faceted search and other search paradigms Other than the faceted search, there are three other widely used search paradigms:

After the implementation few part of system, it given some theoretical as well estimated results. With the comparison of above given three system how FACET is better, the below graph as well table shown in deafly. The below table shows the time required for searching in specific query with proposed as well as existing approach.

**Table 1: Time required in milliseconds for document retrieving to each approach**

Method	5 d	10d	15d	20d
Keyword Search	246	488	723	975
Form Base Search	310	602	923	1178
Directory Search	840	1520	2310	3125
Facet Search	180	352	533	701

**VII. CONCLUSION**

Faceted search is a technique of accessing a large collection of information that is represented by a faceted taxonomy. It enables users to select facets and facet terms to refine the search results in an iterative way. Extensive research has been done in this domain during the last decade. This approach summarized the published literatures, and proposed a faceted taxonomy of research work on faceted search. On the foundation of the taxonomy, three types of facet models, which are based on the set theory, fsystem's key technologies in the framework, namely facet terms extraction.

**VIII. FUTURE WORK**

For the future environment system can focus on personalize search on user feedback sessions as well as recommendation base on user point of interest with database security is the interesting part of system.



## REFERENCES

- [1] D. Dash, J. Rao, N. Megiddo, A. Ailamaki, and G. Lohman. Dynamic faceted search for discovery-driven analysis. In Proceedings of CIKM '08, pages 3{12, 2008.
- [2] Z. Dou, S. Hu, Y. Luo, R. Song, and J.-R. Wen. Finding dimensions for queries. In Proceedings of CIKM '11 , pages 1311{1320, 2011.
- [3]. Haveliwala, T.H., Topic-sensitive PageRank: a Context-sensitive Ranking Algorithm for Web Search. IEEE Transactions on Knowledge and Data Engineering, 2003. 15(4): p. 784-796.
- [4]. Liu, T.-Y., Learning to Rank for Information Retrieval. Foundations and Trends in Information Retrieval, 2009. 3(3): p. 225-331.
- [5]. Ruthven, I. and Lalmas, M., A survey on the use of relevance feedback for information access systems. Knowl. Eng. Rev., 2003. 18(2): p. 95-145.
- [6]. Stoica, E., Hearst, M.A., and Richardson, M., Automating Creation of Hierarchical Faceted Metadata Structures. Proceedings of the Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics (HLT-NAACL), 2007: p. 244-251.
- [7]. Anick, P.G. and Tipirneni, S., The paraphrases search assistant: terminological feedback for iterative information seeking, in Proceedings of the 22nd annual international ACM SIGIR conference on Research and development in information retrieval. 1999, ACM: Berkeley, California, United States. p. 153-159.
- [8]. Ling, X., et al., Mining multi-faceted overviews of arbitrary topics in a text collection, in Proceeding of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining. 2008, ACM: Las Vegas, Nevada, USA. p. 497-505.
- [9]. Dakka, W. and Ipeirotis, P.G., Automatic extraction of useful facet hierarchies from text databases, in 2008 IEEE 24th International Conference on Data Engineering. 2008. p. 466-475,1631.
- [10]. LingPipe Home. Available from: <http://www.alias-i.com/lingpipe/>.
- [11]. Term Extraction Web search – YDN. Available from:  
<http://developer.yahoo.com/search/content/V1/termExtraction.html>.
- [12]. Taxonomy Warehouse. Available from: <http://www.taxonomywarehouse.com/>.
- [13]. Roy, S.B., et al., DynaCet: Building Dynamic Faceted Search Systems over Databases, in 2009 IEEE 25th International Conference on Data Engineering(ICDE), Vols 1-3. 2009. p. 1463-1466, 1768.
- [14]. Zhao, B., et al., TEXplorer: keyword-based object search and exploration in multidimensional text databases, in Proceedings of the 20th ACM Conference on Information and Knowledge Management (CIKM). 2011, ACM: Glasgow, Scotland, UK. p. 1709-1718.
- [15]. Dou, Z., et al., Finding dimensions for queries, in Proceedings of the 20th ACM Conference on Information and Knowledge Management (CIKM). 2011, ACM: Glasgow, Scotland, UK. p.1311-1320.
- [16]. Chen, J. and Li, Q., Concept Hierarchy Construction by Combining Spectral Clustering and Subsumption Estimation, in Web Information Systems – WISE 2006, K. Aberer, et al., Editors.2006, Springer Berlin / Heidelberg. p. 199-209.
- [17]. Oren, E., Delbru, R., and Decker, S., Extending faceted navigation for RDF data, in Proceedings of the 5th International semantic Web Conference (ISWC). 2006. p. 559-572, 1001.