

## A Survey on approaches of Web Mining in Varied Areas

A.Sangeetha<sup>1</sup>, C.Nalini<sup>2</sup>

<sup>1</sup>Ph.D Scholar, Department of Computer Science & Engg., Bharat University

<sup>2</sup>Professor, Department of Computer Science & Engg., Bharat University

---

**Abstract:** There has been lot of research in recent years for efficient web searching. Several papers have proposed algorithm for user feedback sessions, to evaluate the performance of inferring user search goals. When the information is retrieved, user clicks on a particular URL. Based on the click rate, ranking will be done automatically, clustering the feedback sessions. Web search engines have made enormous contributions to the web and society. They make finding information on the web quick and easy. However, they are far from optimal. A major deficiency of generic search engines is that they follow the “one size fits all” model and are not adaptable to individual users.

---

### I. Introduction

The information retrieval goal is to find the documents that are most relevant to a certain Query. The problem of information retrieval is to find the documents that are relevant to an information need from a large document. It deals with notions of Collection of documents, Query (User’s information need), Notion of Relevancy. The types of information’s are text, audio, video, xml structured and documents, source code, application and web services. The types of information needs are Retrospective, Prospective (Filtering). Retrospective means “searching the past”. The different queries are posed against a static collection. Prospective means “Searching the future”. The static queries are posted against a dynamic collection. It is time dependent. The components in information retrieval are user, process, and collection. User- What computer cares about? Process and collection tends to what we care about. The information retrieval cycle consists of five phases. Source selection, query formulation, search, selection and result. The search process consists of Index and document collection. The indexing is a Black box function; its process is not visible.

The main tasks of information retrieval are indexing the documents, process the query, evaluate similarity and find ranking and display the results. The documents are searching that are most closely matching the query. The indexing consists of stop word removal and stemming and inverted index. The removal of stop word usually improves the effectiveness of information retrieval. The lists of stop words are about, afterwards, according, almost, above etc [12].

The stemming is based on suffix stripping. The reason for stemming is that the words that have similar meaning to each other. The stemming removes the some ending of words. E.g.: include, including, includes, included. A porter algorithm is used for suffix striping. The results of indexing are based on some set of weighted keywords. The results of indexing are in the form of [10]:

$$D1 = \{(t1, w1), (t2, w2), \dots\}. \quad (1)$$

Inverted file is used for retrieving the information for higher frequency.

### Problems in Information Retrieval

- How we represent the documents with selected keywords?
- How document and query representations are compared to calculate the weight?
- Mismatching of vocabularies.
- Ambiguous query.
- Depicting of content may be incomplete and inadequate.

The effectiveness of information retrieval can be improved based on keywords. The keywords cover only the part of contents.[13] User can identify the relevant/irrelevant documents based on the weight of the words. We need to be interacting with user and getting the user feedback. The evaluation is based on recall and precision. The more information retrieval process available is open source IR tool kits.

### II. Web Mining

The World Wide Web has been dramatically increased due to the usage of internet. The web acts as a medium where large amount of information can be obtained at low cost. The information available in the web is not only useful to individual user and also helpful to all business organization, hospitals, and some research areas. The information available in the online is unstructured data because of development technologies. Web mining can

be defined as the discovery and analysis of useful information from the World Wide Web data. [14] It is one of the data mining techniques to automatically extract the information from web documents. The three issues in the WWW are web content mining, web structured mining, web usage mining. Web structure mining involves web structure documents and links. Web content mining involves text and document and structures. Web usage mining includes data from user registration and user transaction. WWW provides a rich set of data for data mining. The web is dynamic and very high dimensionality. It is very helpful to generate a new page, lot of pages are added, removed and updated anytime. Data sets available in the web can be very large and occupy ten to hundreds of terabytes, need a large farm of servers. A web page contain three forms of data, structured, unstructured and semi structured data. A number of algorithms are available to make a structured data, one such algorithm is a fuzzy self constructing. An unstructured data can be analyzed using term frequency, document frequency, document length, text proximity.

We have to improve searching in the web by adding structured documents. Using clustering techniques we have to restructure the web information. We provide a hierarchical classification of documents using web directories Eg: Google. While increasing the annual band width in ten times its average is increasing three times, because of that the traffic management is important in web mining.

### **A. Related Works**

In recent years, many works have been done to infer the so called user goals or intents of a query. But in fact, their works belong to query classification. Some works analyze the search results returned by the search engine directly to exploit different query aspects. However, query aspects with no efficiency in have limitations to improve search engine relevance. Some works take personalization through two broader categories namely i)click based methods and ii)profile based methods. The click based methods generate search engine results through clicking of particular link . The most efficient click based methods are ‘classified average precision’ and ‘fuzzy self-constructing method ‘But this strategy works only on repeated queries from the same user. But are not applicable on multiple user queries. On the other hand, profile based methods provide results comparing user profile and user query. It works for hierarchy of user profiles and generate better results .Content based ranking is done is proposed to rank the search engine results by analyzing content and keywords [4]. By analyzing the content and keywords, term frequency is calculated. Term frequency is the number of times a term or a document appears in a page. This determines the total relevancies of a link in a page. This ranking in personalized framework reduces complexity of users and provide better results satisfying all the users.. One application of user search goals is restructuring web search results. There are also some related works focusing on organizing the search.

### **B. RANKING**

Every result page keywords and content words are pre-processed and compared against the dictionary. If a match is found against the content or keyword and the root word then particular weight is awarded to each word. Finally, the total relevancy of the particular link against user request is determined through the term frequency. The page which contains total relevancy value nearest to 1 are ranked as first page and 0 are ranked as last page.

III. Literature Review

Table 1.1 shows that various techniques that are used for retrieve the information from web.

NO	AUTHOR NAME	TITLE	YEAR	DESCRIPTION	ADVANTAGE:	DISADVANTAGE
1	Kiana Calagari, and Ahmed Abdelsadek	Cloud based multimedia content protection system	2016	for large-scale multimedia content protection. two novel components: (i) to assign signatures of 3-D videos, and (ii) distributed matching engines for multimedia objects..	The distributed matching engine leads to high scalability and it is designed to support different multimedia objects	Cost is high.
2	Lei Zhang ,Yongdong Zhang	Full space local topology extraction for cross model retrieval	2015	A novel hashing method based on the withdrawal of the common manifold structure is proposed. This is shared among different feature spaces. two possibilities of local topology information are exploited in this method. Local angles are incorporated within the extraction of local topology of each feature space, which leads to a common intermediate subspace. The proposed method is termed as a full-space local topology extraction for hashing.	A novel cross model hashing method by extract the local data topology within the value space of each modal and between value spaces of different modals	When cross model appear timing is bulky between text and images.
3	Yi-Jie Lu, Linjun Yang	Mining latent attributes from click-through logs for image recognition	2015	Attribute-based image model which it represents an image by prognostic it into a space spanned by attributes, has attracted developing attention from both computer vision and multimedia communities for the compactness and potential to bridge the semantic gap. Different works focus on learning attribute models and utilizing them in image recognition and retrieval, few touch on the problem of how to effectively construct a vocabulary of attributes, which has the essential part of effective attribute-based representatives. the attribute vocabulary by human experts or through existing ontology, which is often limited in coverage of general concept space. The mining of latent topics from the click log is formulated as a matrix factorization problem, and further better by weighted terms-based matrix factorization to address the extreme sparsity of the click-through matrix.	Two qualitative results of the mined LTA and quantitative results on the standard image recognition benchmark demonstrate the mined LTA's effectiveness	The severe accuracy loss is cause by unacceptably defining the attribute terms.
4	Pai peng liden shou	KISS:Knowing camera prototype system for recognizing and annotating places of interest	2016	Knowing camera prototype SyStem (KISS) for places-of-interest (POI) identification and gloss for smartphone photos in the real time systems with the availability of online geotagged images for POIs in knowledge base. "Spatial+Visual" (S+V) framework which consists of a two probabilistic field-of-view (pFOV) model in the spatial phase and sparse coding is proposed. They put a forward an offline Collaborative Salient Area (COSTAR) mining algorithm to identify common visual features (called Costars) among the noisy photos in geotagged on each POI, estimate to clean the geotagged image are stores in the database. The mining result can be utilized to annotate the region-of-interest on the query image during the online query processing. This mining procedure also developed the efficiency and truth of the S+V Framework it is used to extend the pFOV model into a Bayesian FOV(bFOV) model which increased the spatial recognition accuracy by more than 30 percent and also further alleviates visual computation.	There exist numerous study mechanism trying to recognize or retrieve images by chart similarity.	Spatial and visual method cannot address our problem because a photo visually similar to a query may not be properly tagged by a point of interest even if each photo is tagged correct and visually similar photos might be tagged by different point of interest.
5	Richang	Unified photo	2016	Photo enhancement process used to	A tag-wise regularized	In this tag

	Hong, Lumng Zhang	enhancement by discovering aesthetic communities from flickr		Increasing the aesthetic appeal of a photo consists of changing the photo aspect ratio and spatial recombination. They focus on photos from the image hosting site Flickr, which has 87 million users and more than 3.5 million photos are used. First, a tagwise regularized topic model to describe the aesthetic topic of each Flickr user, coherent and interpretable topics are discovered by leveraging both the visual features and tags of photos. Next, a graph is constructed to describe the similarities in aesthetic topics used in users. Notice that densely connected users have similar aesthetic topics, which are categorized into different communities by a dense subgraph mining algorithm. A probabilistic model is exploited on the way to enhance the aesthetic charm of a test photo by leveraging the photographic experiences of Flickr users from the corresponding communities of that photo.	topic model is developed to describe the aesthetic topic of each Flickr user.	based it is difficult to develop a new software like flicker.
6	Chunjie Zhang, Qingming Huang	Contextual exemplar classifier based image representation for classification	2016	A novel contextual classifier based method for image representation and classification tasks is proposed. Each exemplar classifier is trained to separate one training image from other images of different classes. We partition each image into a number of regions and use the responses of these exemplar classifiers as the image region's representation. The contextual relationship is then modeled using mixture Dirichlet distributions. A bi-layer model is used to predict image classes with L2 constraints. Experimental results on the Natural Scene, Caltech-101/256, Flower-17/102 and SUN-397 datasets able to outperform the state-of-the-art local feature based methods image classification.	Exemplar classifiers were taught to split every training image from other images of different classes	The Speed Up Robust Features descriptor which was scale and rotation invariant and was also robust to noise.
7	Dan Lu, Xiaoxiao Liu,	Tag based image search by social re ranking	2016	Social media sharing websites like Flickr used to users images with free tags, which significantly contribute to the enlargement of the web image retrieval and organization. Tag-based image search is used to find images contributed by social reranking. A social re-ranking system is a tag-based image retrieval with the consideration of an image's relevance and diversity is proposed. The re-ranking images according to their visual information, semantic information, and social clues the on line images by inter-user re-ranking. Users that have higher contribution to the given query rank higher. These selected images make up the final retrieved results. We have to build an inverted index structure for the social image dataset to accelerate the searching process.	Images are ranked due to social ranking method.	The top ranked images determined by VR, RR, and CRR, are all suffered from the lack of diversity.
8	Weiming Zhang, Hui Wang,	Reversible data hiding in encrypted images by reversible image transformation	2016	Reversible data hiding in encrypted images (RDH-EI) attracts more Attention for reversible data hiding for encrypted image base on reversible image transformation. Different from all previous encryption-based works in which the ciphertexts may attract the annotation of the curious cloud, reversible image transformation based framework allows the user to transform the original image into the content of another target image with the same size. The transformed image, that looks like the target image, is used as encrypted	Development is extending the traditional upturn to the progressive-based recovery.	Decrypting the clear encrypted image, the generated image has a reduced quality

				image and is outside to the cloud. Therefore, the cloud server can easily incorporated data into the "encrypted image" by any reversible data hiding methods for plaintext images. A client-free scheme for reversible image transformation can be realized, that is, the data-incorporated process execute by the cloud server is irrelevant with the processes of both encryption and decryption. Two reversible data hiding kinds , including traditional reversible data hiding and unified embedding and scramble scheme, are adopted to incorporate watermark in the encrypted image, which can satisfy different needs on image quality and large embedding capacity.		
9	Wei Zhang, Chong-Wah Ngo,	Hyperlink-Aware Object Retrieval	2016	Object retrieval involve small objects on disorderly backgrounds, where the similarity between the querying object and a relevant image can be heavily affected by the background. To address this problem, we propose an efficient object retrieval technique by hyperlinking the visual entities among the reference data set. In particular, a two-step framework is proposed: subimage-level hyperlinking and hyperlink-aware reranking. For hyperlinking, we propose a scalable object mining technique using Thread-of-Features, which is designed for mining subimage-level objects. For reranking, the initial search results are reranked with a hyperlink-aware transition matrix encoding subimage-level connectivity.	Dynamic dataset requires increasing inform hyperlinks, which are not trivial issue and worth more investigation.	Online response time be dangerous for actual applications, we further match up to the efficiency for many retrieval methods.
10	Hanli wang Bo Xiao	A cloud based heterogeneous computing framework for large scale image retrieval	2015	cloud-based heterogeneous computing framework (CHCF), is a set of tools and techniques for collection ,optimization, and execution of multimedia mining applications on different systems. With the aid of the compiler and the utility library provided by cloud based heterogeneous users are able to develop multimedia mining applications .The framework employs a number of techniques, including adaptive data partitioning, knowledge-based hierarchical scheduling, and performance estimation, is proposed to achieve high computing performance. The most important multimedia mining applications, large-scale image retrieval is investigate based on the proposed cloud based heterogeneous framework . The scalability, computing performance, and programmability of cloud based heterogeneous framework are studied for large-scale image retrieval by case studies .	The experimental results exhibit that cloud based heterogeneous framework can achieve good scalability and improvements in image retrieval.	The system needed to find illegal copies of multimedia objects of are complex and large

#### IV. Proposed Work

Web users typically submit very short queries to search engines, the very small term overlap between queries cannot accurately estimate their relatedness. Given this problem, the technique to find semantically related queries (though probably dissimilar in their terms) is becoming an increasingly important research topic that attracts considerable attention.

After the survey and research, it has been found that the need of having a search engine procedure or any searching technique which gives more refined and accurate search results in any of the user defined context. As the various search engines currently present in the market may or may not give the relevant or related search results. So to fill the gap between the output of a search engine from related search results to more related and relevant search results, a technique is required.

The architecture of my proposed research work is represented by a diagram. The implementation has five modules

1. User Profile and Ontology Construction
2. Query mapping and search results
3. Content and keyword extraction
4. Ranking
5. Improved Search Results.

## V. Conclusion

In this paper, first we have mainly focused on the web mining types- Web content mining, web structure mining and web usage mining. After that, we have introduced the web mining techniques in the area of the Web which requires the different goals and also it is useful to develop different business application. Ecommerce is one of the example of this personalization technique which depend on the how well the site owners understood the user's behavior and their needs. Web usage mining is useful for the pattern matching, site reorganization, product/site recommendation etc. Future efforts, investigating architectures and algorithms that can exploit and enable a more effective integration and mining of content, usage, and structure data from different sources promise to lead to the next generation of intelligent Web applications.

## References

- [1]. An effective approach for increasing the efficiency of inferring user search goals with feedback sessions. B. Saranya, G. Sangeetha, Valliammai Engineering College, Chennai.
- [2]. How to Use Search Engine Optimization Techniques to Increase Website Visibility JOHN B. KILLORAN, IEEE TRANSACTIONS ON PROFESSIONAL COMMUNICATION, VOL. 56, NO. 1, MARCH 2013
- [3]. WebCap: Inferring the user's Interests based on a Real-Time Implicit Feedback. Nesrine zemrili, Information system department, 978-1-4673-2430-4/12-2012.
- [4]. A Collaborative Decentralized Approach to Web Search Athanasios Papagelis and Christos Zaroliagis, Member, IEEE IEEE TRANSACTIONS ON SYSTEMS, MAN, AND CYBERNETICS—PART A: SYSTEMS AND HUMANS, VOL. 42, NO. 5, SEPTEMBER 2012.
- [5]. A Web Search Engine-Based Approach to Measure Semantic Similarity between Words, Danushka Bollegala, Yutaka Matsuo, and Mitsuru Ishizuka, Member, IEEE IEEE TRANSACTIONS ON KNOWLEDGE AND DATA ENGINEERING, VOL. 23, NO. 7, JULY 2011
- [6]. Correspondence Falcons Concept Search: A Practical Search Engine for Web Ontology Yuzhong Qu and Gong Cheng IEEE TRANSACTIONS ON SYSTEMS, MAN, AND CYBERNETICS—PART A: SYSTEMS AND HUMANS, VOL. 41, NO. 5, MAY 2011.
- [7]. One Size Does Not Fit All: Towards User & Query Dependent Ranking For Web Databases Aditya Telang, Chengkai Li, Sharma Chakravarthy Department of Computer Science and Engineering, University of Texas at Arlington July 16, 2009
- [8]. Long-Term Cross-Session Relevance Feedback Using Virtual Features Peng-Yeng Yin, Bir Bhanu, Fellow, IEEE, Kuang-Cheng Chang, and Anlei Dong, IEEE TRANSACTIONS ON KNOWLEDGE AND DATA ENGINEERING, VOL. 20, NO. 3, MARCH 2008.
- [9]. Automated Ranking of Database Query Results, Sanjay Agrawal, Surajit Chaudhuri, Gautam Das, Microsoft Research., Aristides Gionis Computer Science Dept, Stanford University, Proceedings of the 2003 CIDR Conference
- [10]. An Efficient k-Means Clustering Algorithm: Analysis and Implementation Tapas Kanungo, Senior Member, IEEE, David M. Mount, Member, IEEE, Nathan S. Netanyahu, Member, IEEE, Christine D. Piatko, Ruth Silverman, and Angela Y. Wu, Senior Member, IEEE, IEEE TRANSACTIONS ON PATTERN ANALYSIS AND MACHINE INTELLIGENCE, VOL. 24, NO. 7, JULY 2002
- [11]. A New Algorithm for Inferring User Search Goals with Feedback Sessions Zheng Lu, Student Member, IEEE, Hongyuan Zha, Xiaokang Yang, Senior Member, IEEE, Weiyao Lin, Member, IEEE, and Zhaohui Zheng
- [12]. [Online]. Available: Introduction to Information Retrieval, Jian-Yun Nie University of Montreal Canada.
- [13]. An effective approach for increasing the efficiency of inferring user search goals with feedback sessions. B. Saranya, G. Sangeetha, Valliammai Engineering College, Chennai.
- [14]. Improved Algorithm For Inferring User Search Goals With Feedback Sessions" in International Journal of Research in Computer Applications and Robotics, A. Sangeetha, C. Nalini, Bharath University Department of Computer Science and Engineering, Bharath University, Tamil Nadu, India
- [15]. Retrieving Relevant Links from the Web Documents through Web Content Outlier Mining From Web Clusters, Volume 5, Issue 2, February 2015 ISSN: 2277 128X International Journal of Advanced Research in Computer Science and Software Engineering A. Sangeetha, T. Nalini Department of Computer Science and Engineering, Bharath University, Tamil Nadu, India