# Opinion Search and Retrieval from WWW

## Dr. A. Padmapriya[1] ,S. Maheswaran[2]

*[1]Department of Computer Science and Engineering, AlagappaUniversity,Karaikudi-Tamilnadu,India*
*[2]Department of Computer Science and Engineering, AlagappaUniversity,Karaikudi-Tamilnadu,India*

**Abstract :** *Opinion retrieval has established itself as an important part of search engines ratings, opinion trends and representative opinions enrich the search experience ofusers when combined with traditional document retrieval by revealing more insights about a subject.In the past years we have witnessed Sentiment Analysis and OpinionMining becoming increasingly popular topics in InformationRetrieval and Web data analysis.With the rapid growth of the user-generated content on the Web.Opinion retrieval is a document retrieving and ranking process, a relevant document must be relevant to the query and contain opinions toward the query. Opinion polarity classification is an extension of opinion retrieval; it classifies the retrieved document as positive, negative or mixed, according to the overall polarity of the query relevant opinions in the document. In this study, we review the development of opinion search and retrieval during the last years, and also discuss the evolution of a relatively newresearch directionand we try to layout the futureresearch directions in the field.*

**Key words -***Opinion mining, Opinion Retrieval, Opinion Identification, Text Mining,*

## I.    Introduction

Since the World Wide Webfirst appeared two decades ago, it has changed the waywe manage and interact with information. It has now become possible to gather the information of our preference from multiple specialized sources and read it straightfrom our computer screen. But even more importantly, it has changed the way we share information. Today people not only comment on the existing information, bookmark pages, and provide ratings, but they also share their ideas, newsand knowledge with the community at large.There exist many mediums, where people can express themselves on the web Blogs,wikis, forums and social networksOpinion search as a research area is arelatively new branch of studies. The aim is to enable users to search for opinions on any object[27]. However, the entity "object" is used to point to different concepts including  products, persons, happenings or topics. Therefore opinion search can be helpful for a broad range of application. Extracting information from news articles and other texts is an important application task for natural language processing technology. In the past few years, web documents are receiving great attention as a newmedium that describes individual experiences and opinions. This situation is generating increasing interest in technologies for automaticallyextracting or analyzing personal opinions from web documents such as posts onmessage board and weblogs. Such technologies can be an alternative to traditional questionnaire based social or customer research and would also be Web userswho seek reviews on certain consumer products of their interest.Previous approaches to the task of mining a large-scale document collectionof customer opinions (or reviews) can be classified into two approaches: textclassification and information extraction approaches. In the former researchershave been exploring techniques for classifying documents or passages accordingto semantic/sentiment orientation such as positive vs. negative [1,19,27] . The latter, on the other hand, focuses onthe task of extracting opinions consisting of information about particular aspectsof interest and the corresponding semantic orientation in a structured form fromunstructured text data. In contrast to sentiment classification, opinion extractionin general aims at producing richer information useful for in depth analysis ofopinions, which has recently been taken on by a growing research community [5],[12],[21].In section 2, the background and the motivation of the current study is described. Section 3 presents the classification framework of the existing work on opinion search. Section 4, contains the discussion of the findings of this study and results are discussed at section 5.

## II.    Background

Information available on the web as text format can be broadly classified into two main categories, fact and opinions. Facts are generally objective statements about entities and events. But opinions are subjective statements that reflect people's sentiments or perceptions about the entities and events which is the area of interest for this work. Most of the exiting research like information retrieval Web search, and other text mining and natural language processing tasks on text information processing has been focused on mining and retrieval of factual information only a little work has been done on the processing of opinion until recently[3,22].Though the  number of research interests in this area is growing fast. As a human being, people like to express their own opinion. They are also interested to know about others opinion on anything they are interested, especially whenever they need to make a decision. The technology of opinion mining thus has a tremendous scope for practical application [26]. In order to enhance customer satisfaction and shopping experience, it has become a common practice for online merchant to enable their customer to review or to express opinions on the products that they have purchased[22].This is better than reading a large number of reviews to form a mental picture of the strengths and weaknesses of the product. Using these reviews, product developing companies can gather feedback on their products in a relatively easy and cheap way. These ratings can be used to evaluate the results of opinion mining techniques by comparing the ratings to the results from opining mining. Activity regarding opinion search has been heavily concentrated. Furthermore, manyreviews are very long and have only a few sentences containing opinions on the product. This makes it harder for a potential customer to read them all to make an informed decision on whether to purchase the product or service. If he/she only reads a few reviews, he/she may get a biased view. It is very difficult for a human reader to find relevant sources, extract pertinent sentences, read them, summarize them and organize them into usable forms [3]. An automated opinion mining and summarization system is thus become important [3, 22].

## III.    Classification Framework on Opinion Search

Opinion mining has been studied by a good number of researchers in very recent years. To make the unambiguous research opportunities in this field, We are classify the current literature in the following ways which is shown in Fig.1.After an extensive study on this area of work, We have identified two main research directions  namely.

### 3.1 Sentimentand Subjectivity Classification
This is the area that has been researched the most in academia. It treats sentiment analysis as a text classification problem. Two sub-topics[3]that have been extensively studied are.
(i)Sentence - level sentiment classification
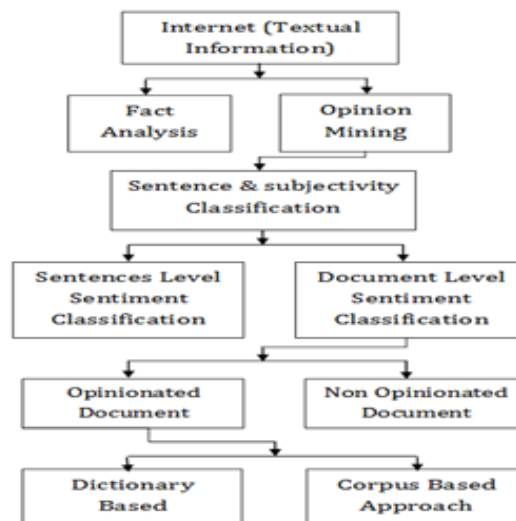
(ii)Document- level sentiment Classification



figure1: classification of opinion mining research.

### 3.1.1Sentence - level sentiment classification

The sentence expresses a single opinion from a single opinion holder. This assumption is only appropriate for simple sentences with a single opinion, e.g., "The picture qualityof this camera is amazing." However, for compound sentences, a single sentence may express more than one opinion. For example, the sentence, "The picture quality of this camera is amazing and so is the battery life, but the viewfinder is too small for such a great camera", expresses both positive and negative opinions[3] (one may say that it has a mixed opinion). For "picture quality" and "battery life", the sentence is positive, but for "viewfinder", it is negative. It is also positive for the camera as a whole.

### *3.1.2 Document level sentiment Classification*

The document level sentiment classification considers as the whole document as the basic information unit. The existing research assumes that the document is known to be opinionated [3].Most existing techniques for document-level sentiment classification are based on supervised learning, although there are also some unsupervised methods [3].

### 3.2 Opinion classification component

This component performs two tasks: (1) classifying each document into one of the two categories, opinionated and not-opinionated, and (2) classifying each opinionated document as expressing a positive, negative or mixed opinion. For both tasks, the system uses supervised learning[3].An opinionated document is a document that contains opinion. The opinion may target at different object. The topic relevant documents of a query are the ideal output of a document retrievalsystem.Such topic relevant documents may or may not contain opinions about the query. But they all contain query related facts.To detect the opinions is the documents, Thus method methods a statistical feature selection process based on the pearson's chi-square text; and the building and applying of a support vector machine(SVM)classifier.The SVM classifier only determines if sentences in a documentsare opinionative or not, but does not know if the opinions areabout the query or not.The SVM classifier has also been shown as one of the the most effective text classification algorithm[28].

The opinionated relevant documents(ORD) of a query are documents that contain query relevant opinions about the query .A list of ORD is the ideal output of an opinion retrieval system. The ORD is a subset of the intersection of the topic relevant documents and the opinionated documents. To find the opinion that are related to the query topic, a novel proximity based method is proposed to evaluate the relevant between a piece of opinion and the query [28].Wei Zhang et al[23] present the **NEAR**operator to classify anopinionative sentence as either relevant or not relevant to thequery.

### 3.3ArchitectureofOpinionSentenceSearch

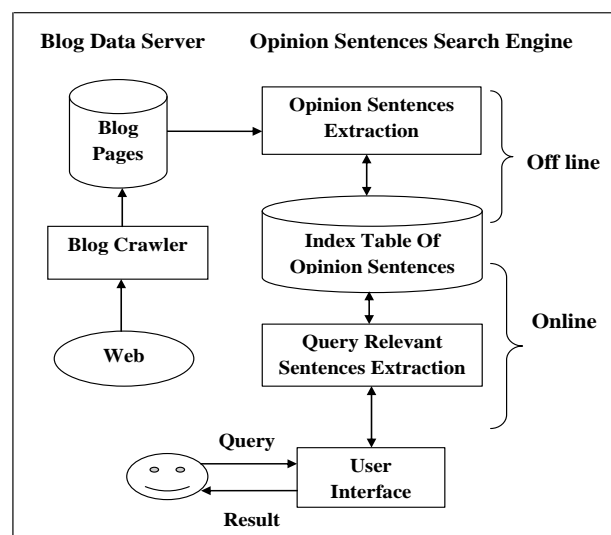The following figure shows the Model of opinion sentence search system in blogspace.



figure2.model of opinion sentence search system.

The blog data server collects blog pages by periodically crawling the web. The opinion sentence search engine, which receives blog pages from the blog data server, consists of two main modules: opinion sentence extraction and query-relevant sentence extraction. The opinion sentence extraction module checks whether each sentence in the crawled blog pages can be considered an opinion. Opinion sentences are extracted and indexed as off-line processing, which, for a practical real-time search, should be as high a proportion of the entire processing as possible. The query-relevant sentence extraction module retrieves opinion sentences relevant to the user's query phrases from the index table of opinion sentences in the blog page server. Since users' queries cannot be predicted, query-relevant sentence extraction has to include on-line processing.

Most document level and sentence level classification methods are based on identification of opinion words or phrases. There are typically two types of approaches:

(i) Dictionary based approach

(ii) Corpus-based approach

Dictionary based approach is based on bootstrapping using a small set of seed opinion words and an online dictionary, e.g., WordNet [11,18]. The strategy is to first collect a small set of opinion words manually with known orientations, and then to grow this set by searching in the WordNet for their synonyms and antonyms. The newly found words are added to the seed list. The next iteration starts. The iterative process stops when no more new words are found. This approach is used in [14,15]. After the process completes, manual inspection can be carried out to remove and/or correct errors. Researchers have also used additional information (e.g., glosses) in WordNet and additional techniques (e.g., machine learning) to generate better lists [2,7,8,10,15]. So far, several opinion word lists have been generated [20,9,14,24,25].The dictionary based approach and the opinion words collected from it have a major shortcoming. The approach is unable to find opinion words with domain specific orientations, which is quite common[3].

Corpus-based approach rely on syntactic or co-occurrence patterns and also a seed list of opinion words to find other opinion words in a large corpus. One of the key ideas is the one proposed by Hazivassiloglou and McKeown [13]. The technique starts with a list of seed opinion adjective words, and uses them and a set of linguistic constraints or conventions on connectives to identify additional adjective opinion words and their orientations[3].It can help find domain specific opinion words and their orientations if a corpus from only the specific domain is used in the discovery process.

## IV. Application Area

Opinion search as a research area is a relatively new branch of studies. The aim is to enable users to search for opinions on any object[27].opinion search can be helpful for a broad range of applications, including review-related websites, forums, blogs, business intelligence, government intelligence and politics, education for e-Learning etc. As most research to date covers opinion search applications in the context of weblogs, review-related websites, Customer reviews on the Internet provide a valuable source of information Reviews available on Amazon.com. provide information on all aspects of products, and on a huge number of alternative products. Activity regarding opinion search has been heavily concentrated on weblogs. They provide a wealth of opinions and information about recent issues regarding a wide range of topics.E-Government refers to the use of information and communications technologies (ICTs) to improve the quality of services and information offered to citizens, to make government more accountable to citizens and advance public sector transparency. Opinion search/retrieval can be used in various fields to meet varied purpose. Binali et al presented some application with some example of current applications [26].

## V. Conclusion and Future work

Opinion search and retrievalhas the potentiality to use from the individual level to organizational level such as companies and government. People and organizations from several domains could benefitted in various way by using the opinion search and retrieval techniques from online customer's feedback. In this paper we have reviewed the current research work in the area of Opinion search and retrieval. We have analyzed several approaches taken by the researchers to extract overall opinion from the unstructured text expressed as opinion even We have critically evaluated the existing work .We strongly believe that this study will help to new researchers to expose cutting edge area of interest in opinion search and retrieval. In the future we plan to apply the decision tree approach to further improve the retrieval effectiveness.

# References

[1]     Ana-Maria Popescu and Oren Etzioni.Extractingproduct features and opinions from reviews..In Proceedings of Human LanguageTechnology Conference and Conference on Empirical Methods in Natural Lan-guage(HLT/EMNLP),2005, 339-346.

[2]     G. Amati. Probabilistic models for informationretrieval based on Divergence from Randomness.PhDthesis, University of Glasgow, 2003.

[3]     Bing Liu,"sentimentanalysi and subjectivity" to appear in Handbook of Natural Language Processing,Second Edition.(editors: N. Indurkhya and F.J. Damerau),2010

[4]     BenHe&JiyinHeIadhOunis"An Effective Statistical Approach to Blog Post OpinionRetrieval"CIKM'08, October 26–30, 2008, Napa Valley, California, USA.

[5]     Bo Pang and Lillian Lee. A sentiment education: Senti-ment analysis using subjectivity summarization based on minimum cuts. InProceedings of the 42nd Annual Meeting of the Association for ComputationalLinguistics(ACL)2004. 271-278.

[6]     Binali,H.,Potdar, V.1 and Chen wul.(2009).A state of the art opinion Mining and Its application domains. IEEE International Conference on Industrial Technology.01/01/2009.

[7]     A. Esuli and F. Sebastiani, "Determining the semantic orientation of terms through glossanalysis," Proceedings of the ACM Conference on Information and Knowledge Management(CIKM), 2005.

[8]     A. Esuli and F. Sebastiani, "Determining term subjectivity and term orientation for Opinionmining,"Proceedings of the European Chapter of the Association for Computational Linguistics(EACL), 2006

[9]     A. Esuli and F. Sebastiani, "SentiWordNet: A publicly available lexical resource for pinionmining," Proceedings of Language Resources and Evaluation(LREC), 2006.

[10]    A. Esuli and F. Sebastiani, "PageRankingWordNetsynsets: An application to opinion mining,"

[11]    C. Fellbaum, ed., Wordnet: An Electronic Lexical Database.MIT Press, 1998.

[12]    Kanayama and Nasukawa,HiroshiKanayama and Tetsuya Nasukawa.Deeper"Sentiment analysis using machine translation technology". In Proc. of the20th International Conference on Computational Linguistics(COLING)2004, 494-500.

[13]    V.Hatzivassiloglou and K. McKeown, "Predicting the semantic orientation of adjectives,"Proceedings of the Joint ACL/EACL Conference, 1997, 174–181.

[14]    M. Hu and B. Liu, "Mining and summarizing customer reviews," Proceedings of the ACMSIGKDDConference on Knowledge Discovery and Data Mining (KDD),2004.168–177.

[15]    J. Kamps, M. Marx, R. J. Mokken and M. de Rijke.UsingWordNet to measure semanticorientation of adjectives.In Proc. of LREC'04,2004,1115-1118.

[16]    S.-M. Kim and E. Hovy, "Determining the sentiment of opinions," Proceedings of theInternational Conference on Computational Linguistics (COLING), 2004.

[17]    Kushal Dave, Steve Lawrence, and David M. Pennock.Miningthe peanut gallery: opinion extraction and semantic classication of productreviews. In Proceedings of the 12th International World Wide Web Conference(WWW), pages 519{528, 2003}.

[18]    A. Lenhart, and S. Fox. Bloggers : a portrait of theInternet's new storytellers.Pew Internet &AmericanLife Project, 2006.

[19]    Minqing Hu and Bing Liu.Mining and summarizing customerreviews.In Proceedings of the Tenth International Conference on KnowledgeDiscovery and Data Mining (KDD), 2004, 168-177.

[20]    I. Ounis, G. Amati, V. Plachouras, B. He,C. Macdonald, and C. Lioma. Terrier: A Higherformance and Scalable Information RetrievalPlatform. In Proceedings of OSIR 2006 Workshop.Proceedings of the Association for Computational Linguistics (ACL), 2007.

[21]    Peter D. Turney. Thumbs up or thumbs down? semanticorien-tation applied to unsupervisedclassication of reviews. In Proceedings of the40th Annual Meeting of the Association for Computational Linguistics (ACL), 2002.417- 424.

[22]    TouhidBhuiyan,YueXu,AudanJosang,"State- of –the-Art Review on Opinion Mining from Online CustomerFeedback.Proceddings of the 9th Asia-Pasific Complex Systems Conference'09.

[23]    Wei Zhang, Clement Yu,WeiyiMeng "Opinion Retrieval from Blogs" CIKM'07, November 2007,6–8,Lisboa, Portugal.

[24]    J. Wiebe, R. F. Bruce, and T. P. O'Hara. "Development and use of a gold standard data set forsubjectivity classifications." Proceedings of the Association for Computational Linguistics ACL),pp. 246–253, 1999.

[25]    P. J. Stone. The General Inquirer: A Computer Approach to Content Analysis. The MIT Press,1966.

[26]    Steven Grijzenhout and ValentinJijkoun and Maarten Marx1 "Opinion mining in Dutch Hansards" the Seventh Framework Programme for Research of the European Commission, under the FET-Open grant agreement FOX, number FP7-ICT-233599.

[27]    Liu,S,Web Data mining :Exploring hyperlinks, contents and usage data.Springer

[28]    Weizhang, Google book of "opinion retrieval and classification in blogs".