



## Opinion Mining of the Movie Blogs based on Supervised Learning Approach

Pranali Yenkar\*

P.G. Student, Department of Computer Engg.  
Datta Meghe College of Engg  
Sector-3,Airoli, Navi Mumbai, . India,400708  
[pranaliyenkar@yahoo.com](mailto:pranaliyenkar@yahoo.com)

Dr. S.D. Sawarkar

Principal, Datta Meghe College of Engg,  
Sector 3-Airoli, Navi Mumbai,  
India, 400708  
[sudhir\\_sawarkar@yahoo.com](mailto:sudhir_sawarkar@yahoo.com)

**Abstract:** An important part of our information-gathering behavior has always been to find out what other people think. With the growing availability and popularity of online review sites and personal blogs, people actively use information technologies to seek out and understand the opinions of others about any product or services. Web Blogs are important new arena for knowledge discovery in open source intelligence gathering. Reviews for products or services in the internet could be in millions which make it difficult to track and understand customer opinions. So in this study, we focuses on the emerging area of research i.e. Opinion Mining that classify the opinions of the users .This paper describes the supervised learning method for the opinion mining of the movie review blogs. The rating and review-summarization system can be extended to other product- review domains easily.

**Keywords :** Blogs, Blog Mining, Opinion Mining, Crawler ,Supervised learning approach

### I. INTRODUCTION

In recent years, blogging has emerged as a popular and important means of communication within Internet users. Blogs, which is short for web-log are a diary format communication channel on the web. They are mostly maintained by individuals, though group blogs are increasingly popular as well. In the blog's entries, the bloggers mention personal or professional information items, references , discuss ideas, interests and opinions. Blogs may be written about any issue,subject or topic, and there is variety of blog types and expressing styles. Due to this properties, analyzing blogs is likely to understand popular culture trend, collect public opinions, and other important information for business. Users also providing their product and services reviews on the Internet platform using blogs. Such reviews are essential and beneficial for customers, merchants, and product manufacturers, marketing personals. According to a survey, 81% of Internet users employ Web to do research about a product they are considering to purchase. Another survey reveals that nearly one-quarter of users (24%) consult consumer reviews prior paying a service and more than three-quarters (79%) of review users report that consumer reviews have a significant influence on their purchase decision. But, as the number of consumer reviews increases rapidly, it becomes tedious for users to obtain a comprehensive view of consumer opinions pertaining to the products of interest through a manual analysis. So Web Blog Mining which is the efficient and effective way of analyzing the sentiments of consumer to specific products becomes desirable and essential which can be used in decision making processes. Opinion mining (also known as sentient extraction, sentiment analysis ) aims to determine the attitude of a speaker or writer with respect to some topic

### II. LITERATURE SURVEY

Etzioni is the first person who introduced the term Web Mining. The web mining subtask and process is described in [1]. The web mining classification like Content mining, Structure mining & Usage mining are described by D. Sravan Kumar and B. Naveena Devi [2]. Web Mining technique has many prominent applications like personalization of the customer experience , analyzing the bidding behavior to determine if a bid is fraudulent, "web wide tracking" which tracks the an individual across all sites he visits. It can provide an understanding of an individual's lifestyle and habits to a level that is unprecedented . The advantages of using web mining in search engines and e-commerce, CRM, customer behavior analysis, cross selling; web site service quality improvement is noticeable and has helped remarkably in Improving the Business Decision Making[3]. Web mining technology also provides the security to the e-commerce web sites. Blogs, a new genre in web 2.0 is a digital diary of web user, which has chronological entries and contains a lot of useful knowledge, thus offers a lot of challenges and opportunities for web mining. In [4],where blog reading behavior of the user is analyzed by blog mining ,they have also defined four kinds of relation between the blogs i.e. citation, Blogroll, Comment, trackback which are called social relation as the relations are publicly observable and therefore involve some degree of social consciousness and manifestation. Opinion mining and detection often called as sentiment analysis, sentiment classification, or sentiment mining automatically identifies emotions in textual data present on the web which is mostly in the unstructured format and extracts sentiment by rating a segment of text as either positive (favorable) or negative (unfavorable).

Sentiment detection helps to enrich business intelligence applications along with how the author of a document perceives a certain product, service, tourism location, or

political party. Mining opinions from the product review is the complicated procedure. First the data needs to be crawled from the web using web crawler then the data needs to be prepared by cleaning it and removing some unwanted tags and non review data and then the data will be mined to summarize the opinion of the users in terms of positive or negative votes or the other way by categorizing as recommended or not-recommended. Some of the past work includes mining reviews of automobiles, banks, travel destinations [6], electronics [5,7] and mobile devices[7]. Pang et al.[8] applied different machine learning approaches such as Naïve Bayes, Support Vector Machines and Maximum Entropy Modeling on movie reviews and obtained considerable results. Turney [6] performs binary classification on product reviews. Like Hatzivassiloglou and McKeon [9] he uses a lexicon containing a set of known sentiment terms which he extends by applying Point wise Mutual Information (PMI) and Latent Semantic Analysis (LSA). The work shows that a simple technique such as PMI is able to outperform the more complicated LSA in such settings.

A more fine-grained approach presented by Pang and Lee determines the exact number of stars provided by the review author. Beineke et al. [10] refine Turney's work [6] by applying a Naïve Bayes model which they train on a labeled and an unlabeled corpus. Like Turney, they use a list of seed terms for the classification of new words, which only contains five positive and negative sentiment terms, as well as a larger list which they assemble from the WordNet [Fellbaum 1998] synonyms of the terms good, best, bad, boring, and dreadful. The authors conclude that their method outperforms previous approaches in regard to classification accuracy and speed of computation. Nicholls and Song [11] examine the impact of different part-of-speech tags by employing a Maximum Entropy classifier. They consider only adverbs, adjectives, verbs and nouns as relevant for sentiment detection and assign these categories different weights. According to their results, adjectives and adverbs are the strongest sentiment conveyors, while verbs and nouns contribute only little. In [12], phrase patterns are used to explain a sentiment classification application which classifies opinions. At the phase of document classification, the tags are added to certain words in the text, and then the tags are matched within a sentence with predefined phrase patterns to get the sentiment orientation of the sentence under study. Next, the sentiment orientation of each sentence is considered and the text is classified according to the sentiment of the most repeated sentiment.

A sentiment miner is described in [13] which uses Natural language processing (NLP) to analyze grammatical sentence structures and phrases and to detect the sentiment of the topic. This method has achieved high quality results (~90% of accuracy) on various datasets including online review articles and the general Web pages and news articles. An application on sentiment classification with review extraction is described in [14]. This uses the sentiment tags and weight approach. Here the review on any particular subject is extracted and a sentiment tag and weight is attached to each expression. Then, it calculates the sentiment indicator of each tag by accumulating the weights of all the expressions corresponding to a tag. Next, it uses a classifier to predict the sentiment label of the text. In this study, the authors use online documents

covering two domains i.e. politics and religion are used to test the performance of the proposed application. The experiments within those domains achieve accuracy between %85 and %95. Opinion mining is studied for the e-learning system in[15].In this study, opinion mining is used to know the users' opinions on the course-ware and teachers of the e-learning system and to help improve the services.

The authors achieve following precisions for these subtasks respectively: %94, %84.2, %80.9 and %92.6. A sentiment mining and retrieval system called Amazing is described in [16]. In this system, the authors incorporate the temporal dimension information into the ranking mechanism, and make use of temporal opinion quality and relevance in ranking review sentences. This study monitors the changing trends of customer reviews in time and visualizes the changing trends of positive and negative opinion respectively. The authors conduct experiments using the customer reviews of four kinds of electronic products including digital cameras, cell phones, laptops, and MP3 players. The evaluation results indicate that the proposed approach achieves a precision of %85 approximately.

The approach proposed in [17] uses a WorldNet, a large lexical database of English, for statistical analysis and movie knowledge. Nouns, verbs, adjectives, and adverbs are grouped into sets of cognitive synonyms. WorldNet is used to generate a keyword list for finding features and opinions. Valid feature- opinion pairs are identified and finally, the sentences are recognized according to the extracted feature-opinion pairs to generate the summary. Experimental results show that this method has an average precision of %65 approximately. An opinion mining application is introduced in [18] which extracts and classifies people's opinions and emotions (or sentiment) from the contents of weblogs about movie reviews to calculate movie scores used unsupervised approach for sentiment mining. In [19] online hotel reviews are mined with supervised machine learning approach using unigram feature to realize polarity classification of documents.

### III. PROBLEM DEFINATION

Internet Web blogs contains lot of text documents that reflects the opinions or sentiments of people about vast entities like book review, software review, political situations, social issues etc. Many people make choices by considering the suggestions, comments and ratings of other people which directly or indirectly influence their purchasing decisions. For example, one likes to buy a CD or DVD that is most recommended by people who use that product. Along with theses websites, a search engine is also an important source for people to search for other people's opinions. If user wants to search anything using search engine, the search engine examines its index and provides a listing of best-matching web pages according to its criteria. However, the semantic orientation of the content, which is very important information in the reviews or opinions, is not provided in the current search engine. For example, Google will return around thousands of hits for the query "3 Idiots review." If search engines can provide statistical summaries from the

opinions point of view, it will be more useful to the user who polls the opinions from the Internet.

A scenario for the aforementioned movie query may yield such report as “There are 10,000 hits, of which 80% are thumbs up and 20% are thumbs down.” This type of service requires the capability of discovering the positive reviews and negative reviews. Thus, there is a need to crawl and process peoples’ opinions, so that it can be used in decision making processes of potential Web review applications. In this study, we propose a web blog mining system that will allow the user to select the movie through the GUI then Crawlers will fetch the movie information from different blogs and then the fetched data is parsed, processed and analyzed to summarize the opinions or sentiments by using supervise learning method. This approach will show Web blog users what other people think about a particular movie by means of graphs or charts. Although this study focuses on movie review, the whole design can be applied to other domains such as restaurant, hotel, etc.

#### IV. PROPOSED METHOD

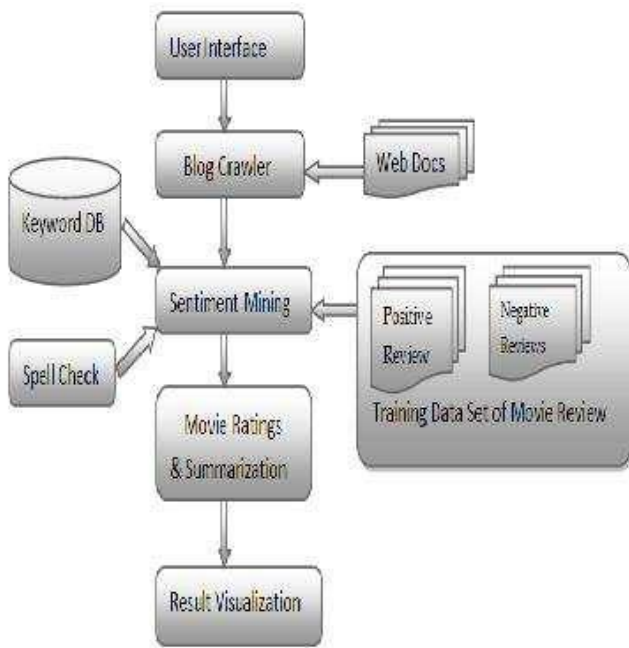


Figure 1: Movie Blog Opinion Mining Architecture

Blog Mining process consists of below mentioned main components;

- User Interface
- Crawling
- Sentiment Mining
- Keyword Database
- Movie Rating and summarization
- Result Visualization

##### A. User Interface:

In this proposed method, a user interface is developed using which a user can access the web blog miner application. User can select the movie name appearing in the drop down list and after pressing the button, the star rating of that particular movie will display along with the additional

information about movie like cast of the movie and release date. User can also view and compare the ratings of more than one movie at a time by selecting the movie names from the list and pressing the Compare button.



Figure 2 :User Interface of the Application

##### B. Crawling:

One of the important components of the proposed system is the Crawler which is also known as Spider. Crawler is software which actually browses the WWW through the lists of the links provided and gathers the specific information from it. Many sites, search engines in particular, use spidering as a means of providing up-to-date data. Web crawlers (or spiders) are mainly used to create a copy of all the visited pages for later processing by a search engine that will index the downloaded pages to provide fast searches. Crawlers can also be used for automating maintenance tasks on a Web site, such as checking links or validating HTML code. Also, crawlers can be used to gather specific types of information from Web pages, such as harvesting e-mail addresses (usually for spam) or gathering text content. Two open source projects, OpenWebSpider and Arachnode can be utilized for crawling the Web blogs and collecting data for sentiment analysis. In this study, the information gathered from the web blog pages after crawling is analyzed and cleaned by removing unwanted no review data like images, links etc. The cleaned data is then provided for the sentiment mining.

##### C. Sentiment Mining:

Sentiment analyzer, opinion miner or sentiment miner is the core component of the opinion mining system. Sentiment analyzer identifies the sentiments or the opinions i.e. agreement or disagreement statements that deal with positive, negative or neutral reviews. Sentiment mining can be done by using two approaches i.e. Supervised Learning Approach and Unsupervised Learning Approach. In our proposed method we are using the Supervised Learning Approach which is better and more accurate to get the result than the



unsupervised learning approach. The crawled movie information will be given to the Sentiment Analyzer as a input.

Sentiment analyzer will then parse the reviews and comments and separate the keywords and then search those keywords in the vast keyword list available in the database. Once the keyword is found then sentiment analyzer will come to know whether the word is positive or negative and its corresponding numeric score. All the keywords relevant to the movie under consideration are searched in the database to find out the accumulative positive and negative scores of the complete sentence and so on for the complete review. The analyzer utilizes the database of the positive, negative, neutral and ignore words and SentiWordNet to obtain the sentiment scores. As mentioned we are using the supervised learning method so the miner is already trained with the number of positive and negative reviews for the movies displayed on the user interface. When the user actually select any movie from the list then Sentiment miner will make use of the training data and determine the sentiment orientation and strength of the sentiment orientation for that particular movie. Spell check facility is also provided to increase the accuracy of the result as peoples' reviews and comments in their blogs may contain spelling errors and these errors will decrease the result of the application

#### D. Keyword Database:

SentiWordNet is an open source lexical resource explicitly developed for supporting sentiment classification and opinion mining applications. SentiWordNet 3.0 is the latest version available for research. We have created and referred a large database comprising of around 21,131 words including positive, negative, neutral, prefix word[20]. Each of these words will have its individual score and category which will be considered for calculating the overall sentiments of the user about the movie.

ID	Score	Category	Synset_term
1	6	prefix	isnt
8	6	ignore	mostly
9	4	ignore	most
583	7	negative	rubbish
584	3	negative	bad
585	4	negative	hate
596	7	negative	annoyed
606	5	negative	awful
711	8	negative	degraded
713	8	negative	dejected
874	10	negative	incomplete
1061	11	negative	notverylike
1062	7	negative	notlike
9976	5	neutral	rapid
9977	4	neutral	rare
10054	4	positive	like
10055	8	positive	awesome
10063	12	positive	captivating
10088	7	positive	lively
10114	10	positive	hilarious
10130	9	positive	splendid
10167	9	positive	blissful

Figure 3 : Keywords Table

#### E. Movie Rating and summarization:

Depending on the positive and negative score range, the rating will be calculated and will be displayed in the form of stars. The star rating of the movie will be proportional to the percentage of positive score. Also the summary of the movie details like Name of the movie, cast etc will display.

#### F. Result Visualization:

Once logged in to the user interface, user can select the movie from the list of the movie names and view the summary of the opinions of the people about the movie in the form of stars. User can also select different movies and view the comparison of ratings in the graphical manner.



Figure 4: Graphical representation of the result

## V. CHALLENGES

Several major challenges apply to Web mining research. First, most Web documents are in HTML (Hypertext Markup Language) format and contain many markup tags, mainly used for formatting. Although Web mining applications must parse HTML documents to deal with these markup tags, the tags can also provide additional information about the document. For example, a bold typeface markup (<b>) may indicate that a term is more important than other terms, which appear in normal typeface. Such formatting cues can be widely used to determine the relevance of terms. Second, traditional IR systems often contain structured and well-written documents (e.g., news articles, research papers, metadata), but this is not the case on the Web particularly for blogs. Web documents are much more diverse in terms of length, structure, and writing style, and many Web pages contain grammatical and spelling errors. Web pages are also diverse in terms of language and subject matter; one can find almost any language and any topic on the Web. In addition, the Web has many different types of content, including: text, image, audio, video, and executable. Numerous formats feature: HTML; Extensible Markup Language (XML); Portable Document Format (PDF); Microsoft Word; Moving Picture Experts group, audio layer3 (mp3); Waveform audio file (wav); RealAudio (ra); and Audio Video Interleaved (avi) animation file, to name just a few.

Web applications have to deal with these different formats and retrieve the desired information. Lastly, the Web

is larger than traditional data sources or document collections by orders of magnitude. The number of indexable Web pages exceeds two billion, and has been estimated to be growing at a rate of roughly one million pages per day. Collecting, indexing, and analyzing these documents present a great challenge. Similarly, the population of Web users is much larger than that of traditional information systems. Collaboration among users is more feasible because of the availability of a large user base, but it can also be more difficult because of the heterogeneity of the user base.

## VI. CONCLUSION

Due the immense growth of the available entertainment, especially movie related websites which has become the major source of the information, the movie goers often overwhelmed with the information. As a consequence, user finds it extremely difficult to obtain any useful comments to make a decision regarding the movie to watch. Hence we have designed a system which actually crawl the movie web blog pages and apply the technique of Sentiment analysis to classify the comments on any particular movie. Supervised learning Approach is proposed for the Sentiment Mining with the extra feature of the Spell Check which further improve the accuracy and performance of the mining.

## VII. REFERENCES

- [1]. Etzioni, O. "The World Wide Web: Quagmire or gold mine", Communication of the ACM, Vol. 39, No. 11, pp. 65-68, 1996.
- [2]. Sravan Kumar, D. and Naveena Devi, B. "Learner's Centric Approach for Web Mining "et al. (IJCSIT) International Journal of Computer Science and Information Technologies, Vol. 1(2), 2010.
- [3]. Mahesh Thylore Ramakrishna, Latha Kolal Gowdar, Malatesh Somashekar Havanur, Banur Puttappa Mallikarjuna Swamy "Web Mining: Key Accomplishments, Applications and Future Directions" 2010 International Conference on Data Storage and Data Engineering
- [4]. Tadanobu Furukawa, Yutaka Matsuo, Ikki Ohmukai, Koki Uchiyama "Analyzing Reading Behavior by Blog Mining" Association for the Advancement of Artificial Intelligence (www.aaai.org)-2007
- [5]. Kushal Dave, Steve Lawrence, David Pennock "Mining the Peanut Gallery: opinion extraction and semantic classification of product reviews" presented at the 12<sup>th</sup> international conference on WWW Budapest, Hungary 2003
- [6]. Peter D. Turney, "Thumbs up or Thumbs Down? Semantic orientation applied to Unsupervised Classification of Reviews" presented at the Association for Computational Linguistics 40<sup>th</sup> Anniversary Meeting, New Brunswick, NJ, 2002
- [7]. Satoshi Morinaga, Kenji Yamanishi, Kenji Tateishi and Toshikazu Fukushima, "Mining product Reputations on the web" presented at the 8<sup>th</sup> ACM SIGKDD international conference on Knowledge discovery and data mining Edmonton, Alberta, Canada, 2002
- [8]. Pang, B., Lee, L., and Vaithyanathan, S. Thumbs up? Sentiment Classification using Machine Learning Techniques. In Proceedings of the Conference on Empirical Methods in Natural Language Processing. Morristown, USA, pp. 79-86, 2002.
- [9]. Hatzivassiloglou, V. and McKeown, K. R. Predicting the Semantic Orientation of Adjectives. In Proceedings of the European Conference of the Association for Computational Linguistics. Morristown, USA, pp. 174-181, 1997.
- [10]. Beineke, P., Hastie, T., and Vaithyanathan, S. The Sentimental Factor: Improving Review Classification via Human-provided Information. In Proceedings of the Annual Meeting on Association for Computational Linguistics. Morristown, USA, pp. 263-269, 2004.
- [11]. Nicholls, C. and Song, F. Improving Sentiment Analysis with Part-of-Speech Weighting. In Proceedings of the International Conference on Machine Learning and Cybernetics. Baoding, China, pp. 1592-1597, 2009.
- [12]. Zhongchao Fei, et al., Sentiment Classification Using Phrase Patterns Proceedings of the Fourth International Conference on Computer and Information Technology (CIT'04), 2004.
- [13]. Jeonghee Yi, et al., Sentiment Mining in Web Fountain, Proceedings of the 21st International Conference on Data Engineering (ICDE 2005), 2005
- [14]. Jian Liu, et al., Super Parsing: Sentiment Classification with Review Extraction, Proceedings of the Fifth International Conference on Computer and Information Technology (CIT'05), 2005.
- [15]. Yun-Qing Xia, et al., The Unified collocation Framework for Opinion Mining,
- [16]. Qingliang Miao, et al., AMAZING: A sentiment mining and retrieval system, Expert Systems with Applications (2008) doi:10.1016/j.eswa.2008.09.035.
- [17]. Li Zhuang, et al., Movie review mining and summarization, Proceedings of the 15th ACM International Conference on Information and Knowledge Management 2006
- [18]. Arzu Baloglu, Mehmet S. Aktas "BlogMiner: Web Blog Mining Application for Classification of Movie Reviews" Fifth International Conference on Internet and Web Applications and Services-2010
- [19]. Han-Xiao Shi, Xiao-Jun Li "A Sentiment Analysis Model For Hotel Reviews Based On Supervised Learning" International Conference on Machine Learning and Cybernetics, Guilin, 10-13 July, 2011
- [20]. <https://github.com/JWHennessey/phpInsight/tree/master/data>