

Analyzing Reader's perception from Top News Blogs using Mining Methodology

Amruta Hingmire¹, Sinju Saliya², Swapnila Mirajkar³

Assistant Professor, Department of Computer Engineering, RSCOE, Tathawade, Pune, India¹

Assistant Professor, Department of Computer Engineering, RSCOE, Tathawade, Pune, India²

Assistant Professor, Department of Computer Engineering, RSCOE, Tathawade, Pune, India³

ABSTRACT: *News Blog* is a user generated website that contains information about current happenings of various fields like world top news, technologies, and sports and so on. It also contains personal opinion of the individual authors. Readers overview or feedback mining matters not only because it contains user's sentiment, but also details about impact of certain decisions, election polls, current trends, popularity and usage of products or service features which could meaningfully change quality of product or service in positive way and helps to capture a more affluent understanding of the mood of an article. In this paper we lay some foundations to aggregate blog articles of a specific area from multiple search services, to analyse the social authorities of articles and blogs, and to monitor the attention articles of the domain obtained over time. This methodology can be combined with additional sentiment analysis methods to create highly automated offline opinion search and engine which can be further used for classification of reviews.

KEYWORDS: News blogs, Review mining, web mining, Indexing, Analysis of opinions.

I. INTRODUCTION

Mining Techniques for On-Line Social Network Analysis

Web mining is an application of data mining, the technique of discovering and extracting useful information from huge data sets or databases. Web mining therefore can be defined as to discover or extract useful information from the web. For on-line social networks analysis, the analysis targets are mainly focused on resources from the web, such as its content, structures and the user behaviours [11]. Thus it can be categorized into three main areas: *content mining*, *structure mining* and *usage mining* [10]. Each one of these areas is associated to the three predominant types of data found in the Web:

Web Content Mining:

Web content mining analyses content on the web, such as text, graphs, graphics, etc. Mining of web content, text or natural language processing are very useful in on-line social network analysis. For example, content mining can group or organize documents on an on-line social networking website, especially articles on blogs or text forums. Article categorization is usually the initial task for many on-line social networks analyses or applications. Web content mining can also be used in on-line social networks analysis to analyse users' reading interests, and determine their favourite content.

News Blog:

News Blog is defined as webpage that contains a series of posts typically characterized by brief texts, providing online commentary, and is presented in reverse chronological order. It's a permanent archive and is searchable. Writing and modifying blog is non technical and so simple that, it encourages everyone for blogging and hence making it popular. The collective community of all blogs is known as the blogosphere. Since all blogs are on the internet by definition, they may be seen as interconnected and socially networked, through blogrolls, comments, and backlinks.

International Journal of Innovative Research in Science, Engineering and Technology

(An ISO 3297: 2007 Certified Organization)

Website: www.ijirset.com

Vol. 6, Issue 5, May 2017

How blog is important for mining?

Many research papers that focus on blogs do not unfurl how comments to the blog posts are taken into consideration. The comments to blog posts vary in terms of length, co references, etc., and thus can be very short answers when the user replies with a short answer or quite long texts when users discuss a topic controversially for instance. By analysing the words and phrases bloggers use while writing their blogs, their mood and state-of mind can be identified and estimated. From our point of view, depending on the type of the blog both the blog posting and the blog comments can be interesting sources for opinion or review mining [4].

Advantages of blog articles:

1. Blogs' use is an interesting marketing technique used by individuals, organizations, and or both.
2. The use of blogs enables consumers and citizens to be better informed about different aspects of life. Since organizations are affected greatly by this, they also tend to create their own blogs in order to affect personal views of people
3. Brand loyalty is an advantage of blogs which benefits organizations since they develop stronger relationship with consumers
4. Blogs' use is a feedback creator, since readers will be more interactive with each other using the "commenting" tool, and people can speak freely about what they think of different aspects
5. Blogs can create the fact of organizations' talent and expertise

RSS feed:

RSS is a way of providing content to the user's browser or desktop in an efficient way. By using RSS feeds, the user can stay updated on the news from electronic paper and other news sources with little extra effort. RSS (Really Simple Syndication) feeds are normally provided in three ways: headlines only, headlines with excerpts and full text feeds.

II. RELATED WORK

Sentiment analysis aka opinion mining is performed at different levels. The analysis in which the complete document is examined is known as document level sentiment analysis. The analysis in which a sentence is examined at a time and then overall analysis is given by summarization of all the sentences is known as sentence level sentiment analysis. n aspect level analysis the feature or entity over which the opinion is given is examined [12]. Aspect level analysis is the most primary level of sentiment analysis. The sentiment words have the most crucial role in identifying the sentiment in a given data Combination of such sentiment words is known as sentiment lexicon [1].The examination of data is done with the help of this sentiment lexicon. One of the challenges faced is analysing sentiments. Another challenge is to identify whether the text is subjective or objective. Subjective sentences convey individual beliefs or feelings, while an objective sentence expresses some true-life information about the world.

Aspect level sentiment analysis was earlier used to extract particular opinion value for a student as proposed in [12]. The drawback in the proposed concept was that only one aspect in each sentence was taken in to account. The second aspect in the sentence under evaluation was ignored. Second flaw was that it didn't take into consideration the fact that a negative sentence can be ended with positive remark with the help of negation and vice versa. This paper removes the flaws in the above concept and negation rules can be used to identify completely positive sentences. The algorithm proposed in the paper takes into account all the aspects given in a particular sentence for finding opinions of individual news reader.

III. PROPOSED SYSTEM

The objective behind such a system is to analyse what are the reviews of people about particular news, about any incident, laws, product launch, famous celebrities; analysis of these results based on their opinions or comments. And generating result showing readers perspective and opinion regarding such events & how such entities are perceived in public.

The system is modelled in two test sets. Firstly we extract the live News articles from the news feeds from various e papers or digital media. Here the subject matter is considered as the news item itself. The live News are extracted and stored locally for analysis and the result is used as a training sample for subjectivity classification. This stage provides a huge amount of valuable information that we are interested to analyse. From fetched textual article, which contains people's opinion, review mining systems analyse:

- What portion contains peoples reviews.

International Journal of Innovative Research in Science, Engineering and Technology

(An ISO 3297: 2007 Certified Organization)

Website: www.ijirset.com

Vol. 6, Issue 5, May 2017

- Who wrote the particular comment?
- What is being commented?
- Generate result for entered keyword.
- Classification of positive and negative feedback/comment.

The system uses lexicon based approach for identifying mood of comments. In this approach, the features in sentences are considered for opinion retrieval. A sentence may contain several features. Different features possibly will have different opinions. For example, the look of Camera is great (positive), but the camera is heavy (negative). All the opinion words in a sentence are identified which can be further used for generating report showing overall analysis.

IV. SYSTEM ARCHITECTURE

A. Opinion Mining

Opinion Mining (OM) is a field of web content mining that aims to find valuable information out of user's opinion. Mining opinions on the web is fairly new subject, & its importance has grown significantly mainly due to fast growth of e-commerce, blogs and forums [4]. More informally, it's about extracting the opinions or sentiments given in a piece of text.

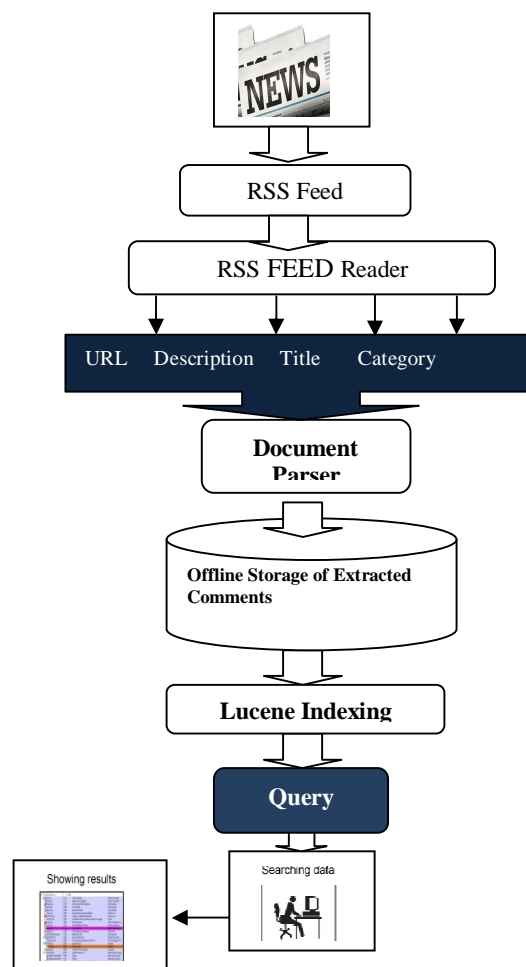


Fig. 1 Working of Opinion mining tool

International Journal of Innovative Research in Science, Engineering and Technology

(An ISO 3297: 2007 Certified Organization)

Website: www.ijirset.com

Vol. 6, Issue 5, May 2017

The result obtained from opinion mining tools will help to identify how a specific entity, a technology, product or any political issues perceived in public.

Item 1:

Title: Sidewalk contract awarded

Description: The city awarded the sidewalk contract to Smith Associates. This hotly contested deal is worth \$1.2 million.

Link: <http://www.gardencitynews.com/contractawards/sidewalk.htm>

Item 2:

Title: Governor to visit

Description: The governor is scheduled to visit the city on July 1st. This is the first visit since the election two years ago. The mayor is planning a big reception.

Link: <http://www.gardencitynews.com/news/2010/06/gov-visit.htm>

Fig 2: Simple RSS channel

B. Data Aggregation Process

For our analyses, we need the URL of each article along with the date of publication, the title and the Textual content and use a XML parser that reads the xml feed from specific news blogs of our interest. Bloggers use RSS (Really Simple Syndicate) feeds to automatically notify their subscribers about new blog posts. This format is very well organized and very frequently used by the bloggers. Our RSS crawler can from time to time retrieve the latest news posts from the news websites using the RSS feeds of the individual sites. As the methodology is intended to monitor a domain over a very long period of time, the crawler is implemented as a permanently running service that regularly queries the search services for the latest articles, and adds these to the data set [6].

C. Cluster analysis

This technique used to group data into sets of elements that are meaningful and/or useful. In this case, clusters represent classes, or conceptually meaningful groups of objects that have familiar characteristics. In other words, clustering is a technique to automatically discover classes. An example of clustering in this context is seen in Information Retrieval. In this case, clustering can be used to group reviews collected from article and are merged together under one group which tells talks about specified topic. This cluster also contains one header file which is associated with the title of new article.

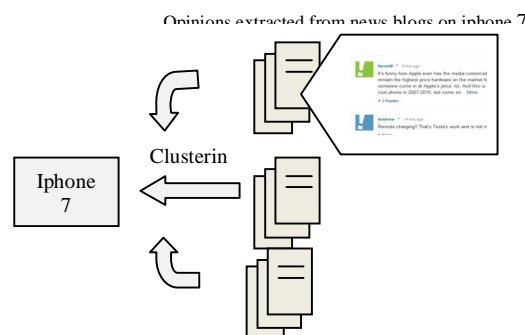


Fig. 3 Reviews are combined from multiple news blog sites which talks about iphone 7 and its feature & are clustered under a name "iphone 7".

International Journal of Innovative Research in Science, Engineering and Technology

(An ISO 3297: 2007 Certified Organization)

Website: www.ijirset.com

Vol. 6, Issue 5, May 2017

Then users only need to inspect the clusters appropriate to their queries, which dramatically decrease the time and the efforts spent in sifting through the lengthy list of documents and prediction of the value of given data at the specific time in future is used for trend prediction. The methods of Web text mining are similar to mining of flat text files to some extent.

Syntax	Description
Token t	Match tokens t
Phrase "cs"	Match tokens in cs in exact order without gaps
cs1?cs2	Match tokens starting with cs1, ending with cs2, with any char between
cs1*cs2	Match tokens starting with cs1, ending with cs2, with any char sequence between
Q1 OR Q2	Match query Q1 or query Q2 (or both)
Q1 AND Q2	Match query Q1 and match query Q2
Q1 NOT Q2	Match query Q1 but not query Q2
+P	Token or phrase P must appear
-P	Token or phrase P must not appear

TABLE 1: Query Language Syntax

However, additional information is conveyed by the tags in Web documents, such as <Title>, <a>, <Heading> and so on which can be exploited to increase the quality of Web text mining.

D. Opinion Lexicon

We have used opinion lexicons as a resource that associates sentiment orientation and words. Their use in opinion mining research stems from the hypothesis that individual words can be considered as a unit of opinion information, and therefore may provide clues to document sentiment and subjectivity [5]. Manually created opinion lexicons were applied to sentiment classification, where a prediction of document polarity is given by counting positive and negative terms. In order to arrive at sensible conclusions, sentiment analysis has to understand context. For example, "fighting" and "disease" is negative in a war context but positive in a medical one. To resolve these confusion domain specific lexicons can be created for each such entity. The query language syntax can also be helpful for finding exact matching token or a keyword.

Opinion orientation can be classified as belonging to opposing positive or negative polarities – positive or negative feedback about a product, favourable or unfavourable opinions on a topic – or ranked according to a spectrum of possible opinions, for example on film reviews with feedback ranging from one to five stars.

International Journal of Innovative Research in Science, Engineering and Technology

(An ISO 3297: 2007 Certified Organization)

Website: www.ijirset.com

Vol. 6, Issue 5, May 2017

E. Sentiment Analysis

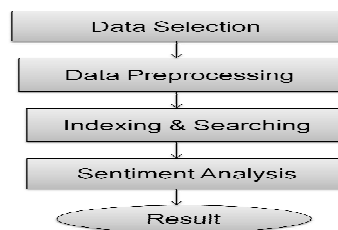


Fig.4: Sentiment Analysis Model

Opinions/comments extracted from blogs need to be categorized based on its polarity i.e. either positive or negative. This task can be carried out by examining the sentiment encoded by text [10]. Data pre-processing steps cleans the uncertain data which is not needed further. Once data is stored locally indexing module will index the documents by using incremental indexing. Sentiment analysis or subjectivity classification can be defined as automated extraction of subjective content from digital text and predicting subjectivity such as positive or negative. It deals with predicting the sense of the post text. For example, if user has written a movie review then sentiment analyser tries to predict whether user is criticizing the movie or he is praising the movie [8].

Example:

Keep it up! Up vote. → **Positive**

The sales will go down → **Negative**

This analysis of sentiment can be performed on several levels of granularity [7] (word/feature, sentences & Documents). Both the document level and the sentence level analysis do not find out what exactly people liked and did not like. Feature level performs finer-grained analysis. In this work this analysis is carried out at the feature & document level.

NLP techniques:

NLP is a way for computers to analyse, understand, and derive meaning from human language in a smart and practical way. By utilizing NLP, developers can categorize and structure knowledge to perform tasks such as automatic summarization, translation, named entity recognition, relationship extraction, sentiment analysis, and topic segmentation. **Open Source NLP Libraries** provide the algorithmic building blocks of NLP in real-world applications. Example of such tool is NTLK [15]

Natural Language Toolkit (NLTK): Python library that provides modules for processing text, classifying, tokenizing, stemming, tagging, parsing, and more. Some example of NLP rules are discussed below.

Negation Rules

A negation word such as no, not and never and also some words that follow patterns such as \stop" + \vb-ing", \quit" + \vb-ing" and \cease" + \vb-ing" change the orientation of opinion words in the following way:

- I. Negation Negative → Positive
- II. Negation Positive → Negative
- III. Negation Neutral → Negative

Some examples for each negation rule:

- I. no problem"
- II. Not good"
- III. Does not work"

TOO Rules

The word \too" is usually used to give a negative connotation if found before adjectives.

The idea is to apply this rule to words which have orientation dependent on context.

- I. The battery life lasts long"
- II. The initialization time takes too long."

Therefore whenever a too word is found before an opinion word, the orientation of the opinion should be negative.

International Journal of Innovative Research in Science, Engineering and Technology

(An ISO 3297: 2007 Certified Organization)

Website: www.ijirset.com

Vol. 6, Issue 5, May 2017

V. EXPERIMENTAL RESULTS

We have analysed number of news blogs from TOI [14] for extraction of opinions. With the help of designed search interface as shown in fig. 5 user can look for reader's reviews or a feedback for particular aspect of topic or he can also search for positive and negative subjectivity of result to directly classify result for sentiment analysis. The same interface can be used for trend detection, popularity check, analysing impact of any legal laws, famous Politician and for any product or service review by analysing blogs written on it.

Companies, organizations and market researchers may benefit from such tools to understand how their product or service is perceived in public, which will help them to predict and increase sales of their product/service.

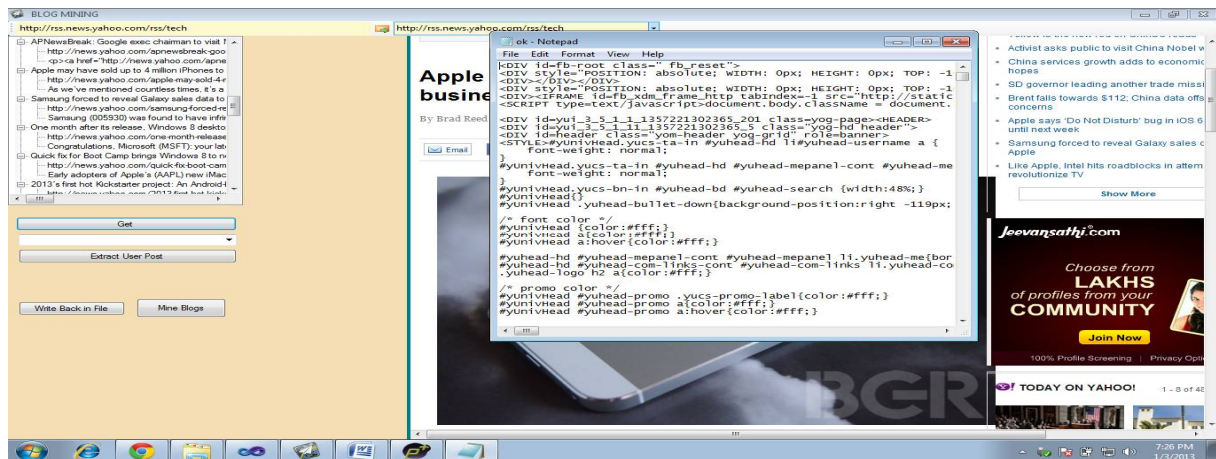


Fig 5: Page source of fetched blog

With the help of above interface news blogs are fetched using RSS feed and reviews are stored which can be used for offline analysis. Multiple blog sites can be accessed for comment extraction.

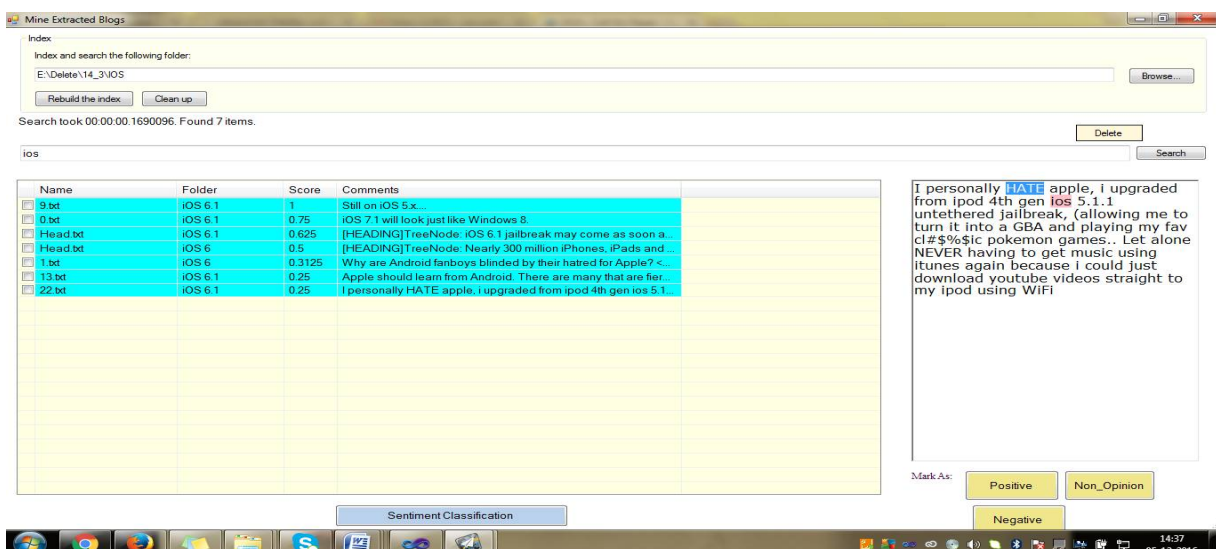


Fig 6: Result set showing opinions containing user entered keyword.

International Journal of Innovative Research in Science, Engineering and Technology

(An ISO 3297: 2007 Certified Organization)

Website: www.ijirset.com

Vol. 6, Issue 5, May 2017

VI. CONCLUSION AND FUTURE WORK

Using review mining tool it is easy to categorize reader's views, ideas and demands. This will be helpful for government, market researchers, developers and industries to analyse their opinion or feedback for decision making. Emoticons, such as :) :-)-and :(, are frequently used online in social media, IM (e.g., Skype), blogs, forums, and other kinds of online social interactions. Because they are commonly used in online communications and they are often direct signals of sentiment, emoticons in text which can be used for understanding reader's opinion or interest. This work can be extended further for such sentiment which will be used to understand opinions in the blogs in more effective way.

REFERENCES

- [1] N. Agarwal, H. Liu, L. Tang, and P. S. Yu. Identifying the influential bloggers in a community. In Proceedings of the First ACM International Conference on Web Search and Data Mining (Video available at:http://videlectures.net/wsdm08_agarwal_iib/), 2008.
- [2] D. Gruhl, R. Guha, R. Kumar, J. Novak, A. Tomlins, The Predictive Power of Online Chatter, in KDD 05 Proceedings of the 11th ACM SIGKDD International Conference on Knowledge Discovery in Data Mining, 78–87, ACM Press, New York, 2005.
- [3] B. Pang, L. Lee, and S. Vaithyanathan, “Thumbs up? Sentiment classification using machine learning techniques,” in Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP), pp. 79–86, 2002.
- [4] Gerald Petzl, Michał Karpowicz, Harald, Andreas Auinger, “Opinion Mining on the Web 2.0 – Characteristics of User Generated Content and Their Impacts” University of Applied Sciences Upper Austria, Austria.
- [5] Bruno Ohana and Brendan Tierney, “Sentiment Classification of Reviews Using SentiWordNet”, Dublin Institute of Technology, School of Computing Kevin St. Dublin 8, Ireland.
- [6] Darko Obradovic, Stephan Baumann, “A Social Network Analysis and Mining Methodology for the Monitoring of Specific Domains in the Blogosphere”, Odense, Denmark.
- [7] Erik Tromp and Adversitement B.V.,” RBEM: A Rule Based Approach to Polarity Detection” Uden, the Netherlands.
- [8] Mukul Joshi and Nikhil Belsare, “BlogHarvest: Blog Mining and Search Framework” International Conference on Management of Data COMAD 2006, Delhi, India, December 14-16, 2006 ©Computer Society of India, 2006.
- [9] Daniel E. O’Leary, “Blog mining-review and extensions: From each according to his opinion”, Marshall School of Business, University of Southern California, Los Angeles, CA 90089-0441, United States.
- [10] Barbara Pobleto (2009) —”Query-Based Data mining for the Web”, TESI DOCTORAL UPF, University Pompeu Fabra.
- [11] I-Hsien Ting “Web Mining Techniques for On-line Social Networks Analysis”, Department of Information Management, National University of Kaohsiung.
- [12] Virmani Deepali, Malhotra Vikrant, Tyagi Ridhi , “Aspect Based Sentiment Analysis to Extract Meticulous Opinion Value” in Proceedings of International Journal of Computer Science and Information Technologies, 2014, pp.3262 – 3266.
- [13] <http://technorati.com/blogs/directory/>
- [14] <http://blog.search.yahoo.com/>
- [15] <http://www.nltk.org/>