

Analyzing and Interpreting Variations of Public Sentiments on Social Networking Site: Tweeter

Mr.Gurunath G.Machhale

Department of Information Technology
Siddhant College of Engineering, Sudumbare,Pune
gurunath.machhale@gmail.com

Prof.Rashmi Deshpande

Department of Information Technology
Siddhant College of Engineering, Sudumbare,Pune
Rashmi2810@gmail.com

Abstract—Due to rapid growth of internet user , Social Networks have become one of the admired communication medium used over internet. Millions of messages are appearing regularly on popular web-sites that provide web services such as Twitter , face book, LinkedIn . Millions of users share their personal opinions or views about on various of issues and discuss several current hot topics on Twitter, making it a important base for tracking and analyzing sentimentation of society. Twitter is a novel micro-blogging platform with more than 23 million unique monthly visitors. On Twitter, any user used to publish a message called as tweet, which is visible publically. Such tracking and analysis can provide critical information for decision making or opinion mining in variety of domains. In this our work, we have moved one step further to interpret sentiment variations. We observed that emerging topics (named foreground topics) within the sentiment variation periods are highly related to the genuine reasons behind the variations. we propose a Latent Dirichlet Allocation (LDA) based model, Foreground and Background LDA (FB-LDA), to distill foreground topics and filter out longstanding background topics. These foreground topics helps to achieve interpretations of the sentiment variations on social networking sites. We select the most representative tweets for foreground topics and develop another generative model called Reason Candidate and Background LDA (RCB-LDA) to rank them with respect to their “popularity” within the variation period.

Sentiment analysis also known as opinion mining plays a crucial role in determining the sentiments involved in various Web content. Analyzing opinions is very important for making decisions. For example, if one wants to buy a new cell

phone, a Web savvy buyer will almost always first check reviews about it in order to make an informed buying decision based on others experiences. Sentiment analysis is currently a very significant trend in the area of natural language processing. Natural language processing involves giving artificial intelligence to computers and is concerned with promoting an understanding of human languages for machines' use. Sentiment analysis extracts opinions, sentiments, and emotions from text and analyses them this information is very useful for governments, businesses and individuals. While this content meant to be helpful in analyzing this bulk of user generated content is difficult and time consuming. So need arises to develop an intelligent system which mine such huge content and classify them into Negative, Positive, Neutral type. Sentiment analysis is the automated mining of opinions, attitudes, emotions from text, database sources through Natural Language Processing (NLP).

Keywords— Sentiment Analysis, Public sentimental Social Sites, Twitter, Emerging topic mining.

I. INTRODUCTION

With the Extensive growth of user generated messages on internet, Social site like Twitter where millions of users used to share their opinion regarding some topic. Figure 1 indicates that web has huge amount of data and social networks has part of that huge data. We can

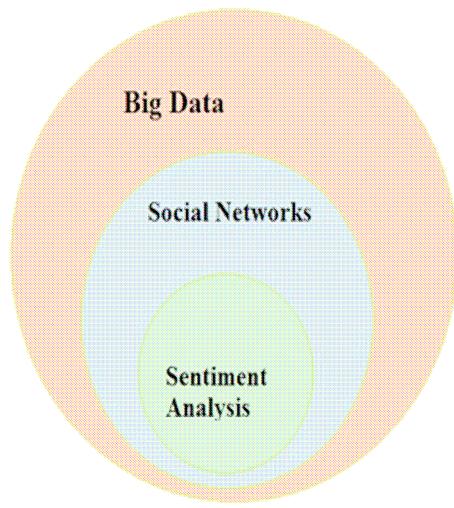


Fig. 1 Visualization of Social Network Analysis

Sentiment analysis on Social Sites data has provided a platform where timely public sentiment can be exposed in an economical and effective way, which is tedious for decision making in various domains. For example, a company could analyze the public sentiment in Tweets to obtain users' feedback towards its products or service; while a political leader can adjust his/her position with respect to the sentiment change of the public, opinion about movies can be most useful for succession of movies Sentiment. Analysis is a technique for extracting sentiment associated with polarities of positivity, negativity and neutrality. It is one of the types of natural language processing in which we can track the mood of the public about a particular entity. Sentiment analysis, which is also called opinion mining, is used for constructing analysis system to collect and examine opinions about the entities on tweets on Twitter. Due to the tremendous of social media services there is great opportunity to understand and analyze the sentiment of the public by analyzing its large-scale data as well as opinion-rich data. Sentiment analysis on tweets can be done by many approaches. Various methods such as machine-learning and lexicon-based approaches have been widely used for sentiment analysis on Twitter like sites. Machine-learning approaches to sentiment analysis need to train the data.

Searching for people's opinions via surveys and polls has been an expensive and time-consuming task. The proliferation of Web 2.0 has changed

the way people express their opinions and feelings. This so called user-generated content posted in blogs, forums, product review sites and social networks is mostly publicly available and easy to obtain. The high value of this content arises from its subjective nature which, in aggregated form, indicates public opinion. It is difficult for humans to read and summarize all relevant documents in terms of the expressed sentiment. Thus, there is a growing need for automated analysis of this kind of data. This is a challenging task with foundations in natural language processing and text mining referred to as sentiment analysis. Many research studies in sentiment analysis are concerned with product reviews from websites like Twitter is a most popular worldwide social website, which provides a micro blogging services and social networking, enable its users to update their status in tweets, follow the people they are interested in(e.g. Sachin Tendulkar) and retweet other's posts and even communicate with them directly. The public Sentimental analysis on Twitter data has provided an economical and effective way to expose timely public sentiments, which is critical for decision making in various domains areas. For instance, a company can study the sentiments of public in Tweets to obtain users 'feedback towards its products. Tweeter is one of the most popular social networking websites, which is drawing more and more attention from researchers from different disciplines. There are several streams of research investigating the role of Twitter. Twitter has attracted attention in both academia and industry for Research Area. Previous research mainly focused on tracking public sentiment.

There have been a large number of research studies and industrial applications in the area of public sentiment tracking and modeling. Previous research like O'Connor [1] focused on tracking public sentiment on Twitter and studying its correlation with consumer confidence and presidential job approval polls. On Twitter, any user can publish a message referred to as tweet, which is visible on the public display.

Similar kinds of studies have been done for investigating the reflection of public sentiments

on oil price indices and stock markets. They reported that events in real life indeed have a significant and immediate effect on the public sentiment on Twitter. One valuable analysis is to find possible reasons behind sentiment variation, which can provide important information for decision-making. E.g. if negative public sentiment towards Barack Obama increases significantly, the White House Administration Office may be eager to know why people have changed their opinion and then react accordingly to reverse this trend. Another example is, Analyzing public opinion variation polling for Exit poll for any Election.

II. LITERATURE SURVEY

Several Researchers carried out research work in Social Network Analysis and sentiment analysis. SA is a text processing technique to derive an opinion or mood intention based on the terms used in a real language sentence. The numbers of researchers have concentrated on generating statistical inference from social network data using sentiment analysis models. Bo Pang and Lillian Lee [2] provided an insight full discussion on sentiment analysis. In this they have considered the ratio in positive words and total words to estimate the opinion.

Today's users can easily obtain information but also they can actively generate content. News reports, BBS, forums, blogs, and etc are the main sources of public opinion information. The text from these sources can contain both facts and opinion which could be extracted using natural language processing mechanisms. Opinions are usually subjective expressions that describe sentiments or feelings of people toward entities and events. It is a sub-discipline of computational linguistics that focuses on extracting opinion of people from the web.

Social media technologies take on many different forms including social network, micro blogging, weblogs, magazines, Internet forums, social blogs, photographs, video, rating and social bookmarking. Micro blogging websites have evolved to become source of various kinds of

information. Due to feature of micros blogs on which people post real time messages about their opinions on a variety of topics, discuss current issues, problems, express positive, negative sentiments for person, events or for products they use in daily life. The manufacturing Companies of such products have started to poll these micro blogs to get a sense of public sentiment for their product. Public and private opinions about variety of subjects are expressed and spread continually via numerous social media. Sentiment analysis is used to determine the attitude of a writer with respect to some topic. The attitude may be his or her judgment, the intended emotional communication or the emotional state of the author when writing. A basic task in sentiment analysis is classifying the polarity of a given text at the word, sentence, document whether the expressed opinion in a word, sentence or in document has sentence feature positive, negative, or neutral. Classification of Sentiments looks, for instance, at emotional states such as 'happy', 'angry', 'sad' and 'neutral'. Sentiment analysis has become popular in judging the opinion of consumers towards various brands [5]. The way in which consumers express their opinion on social networking websites helps to judge this opinion [6]. When it comes to sentiment or opinion or emotion we are not concerned with the topic of the text but the positive or negative opinion it express. People can freely express their opinion in social media as blogs, micro blogs, reviews, forum discussion and social network sites towards any person, events, product, service, news or organization. All these platforms are source of huge amount of valuable information that we are interested to analyze.

Several prior studies have estimated and made use of aggregated text sentiment. The informal study by Lindsay (2008) focuses on lexical induction in building a sentiment classifier for a proprietary dataset of Face book wall posts (a web conversation/micro blog medium broadly similar to Twitter), and demonstrates correlations to several polls conducted during part of the 2008 presidential election. We are unaware of other research validating text analysis against traditional opinion polls, though a number of

companies offer text sentiment analysis basically for this purpose (e.g., Nielsen Buzzmetrics). There are at least several other studies that use time series of either aggregate text sentiment or news, including analyzing stock behavior based on text from blogs, good and bad news (Gilbert and Karahalios 2010), news articles (Koppel and Shtrimberg 2004; Lavrenko et al. 2000) and investor message boards (Antweiler and Frank 2004; Das and Chen 2007). Dodds and Danforth (2009) use an emotion word counting technique for purely exploratory analysis of several corporations.

Twitter is making it a valuable platform for tracing and analyzing public sentiment. It provides information for decision making in various domains with respect to the current issues in society. In this work, we interpret sentiment variations over various topics from society. An recent topics within the sentiment variation periods are related to the genuine reasons behind the variations. Based on this observation, Latent Dirichlet Allocation (LDA) based model, Foreground and Background LDA (FB-LDA), to distill foreground topics. It filters out longstanding *background topics*. Foreground topics can give interpretations about the levels of sentiment variation. This proposed system selects the most representative tweets data for foreground topics and develop model called Reason Candidate and another generative model called Background LDA (RCB-LDA) to rank them with respect to their popularity‘within the variation period. Latent Dirichlet Allocation (LDA) based models to analyze tweets in significant variation periods, and identify possible root cause for the variations. This model can be termed as Foreground and Background LDA (FB-LDA) Mode 1, and it can filter both background topics and extract foreground topics from tweets in the specified variation period, by the use of an auxiliary set of background tweets generated just before the variation. Reason Candidate and Background LDA (RCB-LDA). RCB-LDA first extracts representative tweets for the foreground topics (obtained from FB-LDA) as reason candidates. After that it will associate each remaining tweet in the variation period with one reason candidate and rank the reason candidates

by the number of tweets associated with them. (21-2 night)(80%)

There are many papers, which describe different classification techniques for sentiment analysis. Sentiment classification can be formulated as a supervised problem with two class labels (positive and negative). In (Pang, Lee and Vaithyanathan 2002), the authors apply supervised learning methods such as naïve Bayesian and support vector machines (SVM) to classify movie reviews into two classes. Most unsupervised sentiment classification approaches try to generate a general or domain dependent opinion lexicon for words or opinion phrases. In (Riloff and Wiebe 2003), the authors collected subjectivity clues as a part of their work. The clues were then used in (Wiebe, Wilson and Cardie 2005) to detect semantic orientation. In this paper, a bootstrapping process was proposed, where high precision classifiers use known subjective vocabulary to separate subjective and objective sentences from a no annotated text collection. The aspect extraction method refers to the concept of determining opinion targets and their attributes which are mentioned in a document or a sentence. Many information extraction techniques have been applied so far.(91%)

III. PROPOSED METHODOLOGY

A. General Architecture

Today’s almost all Social Networking sites have been widely used for expressing opinions or emotion in the public domain with help of internet. And Twitter has been the point of attraction to several researchers in important areas. Sentiment analysis over Twitter offers a fast and efficient way to analyze the public sentiment. The main two-fold contributions of this paper are: (1) to the best of our study, our research and our knowledge is the first work that tries to analyze and interpret the public sentiment variations in micro blogging services like twitter. (2) Two novel generative models are developed to solve the reason mining problem. The two proposed models are general: they can be applied to other tasks such as finding topic differences between two sets of documents.

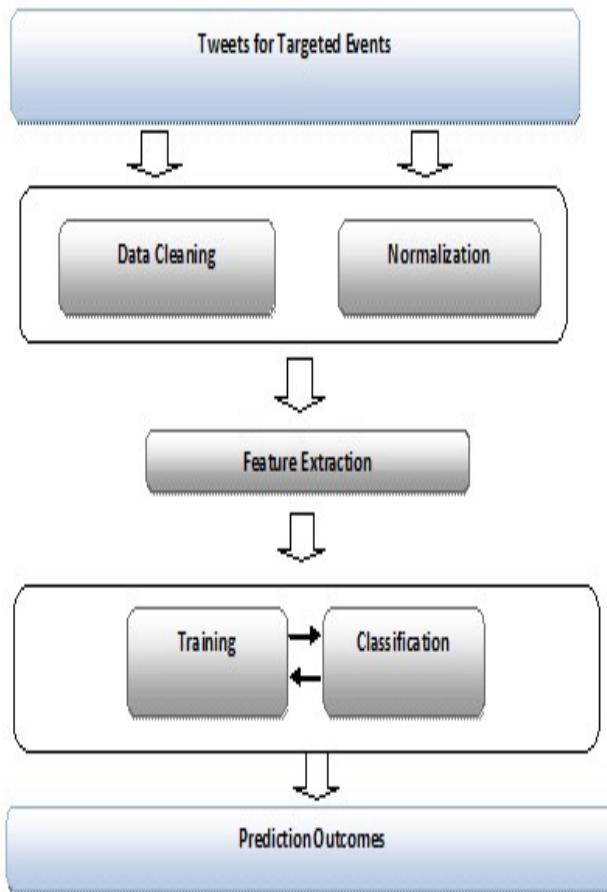


Fig. 2 High Level System Flow

Fig. 2 shows an example of High level system flow. To analyze variations in public sentiments There are two Latent Dirichlet Allocation (LDA) based models: (1) Foregroundand Background LDA (FB-LDA) and (2) Reason Candidate and Background LDA (RCB-LDA). NaïveBayes, SVM, MaxEnt, ANN classifiers with features extracted from Twitter data using feature extraction methods such asUnigram, Bigram and Hybrid (Unigram + Bigrams) for sentiment analysis. In order to remove stop words and to extract features from text, we perform data cleaning and normalization on given set. We extract the target based extended features model [7] by modifying it and twitter user data from the normalized data. vectors are used in part of chunks to train the classifier as a part of incremental training. The sentiment analysis results are incorporated with influence factor of supervised learning to predict the results using prediction model .

B. Proposed Architecture

In our work, we have proposed following three steps for sentiment tracking.

- We extract tweets related to our interested targets (e.g. Arvind Kejriwal, Delhi Election 2015 etc), and preprocess the raw extracted for more Cleaned for sentiment analysis.
- Second, we assign some label so called sentiment label for every individual tweet by combining two state-of-the-art sentiment analysis tools [9], [8].
- finally, depend upon the sentiment labels obtained for each tweet, we identify the sentiment variation for the corresponding targeted issues by using some descriptive statistics.

IV. MODULES

A. Tweets Extraction and Preprocessing: Our First phase starts with extracting tweets lines related to the targeted issue, we go through the whole collected raw dataset and extract all the core lines tweets which contain the keywords of the targeted issues. Compared with regular text documents, tweets are generally somewhat Informal and often written in an adhoc manner like it may contains short forms, some abbreviation. Sentiment analysis can tools applied on raw tweets but often achieve very poor performance in most cases. Hence there is need of preprocessing techniques on tweets are necessary for obtaining satisfactory results on sentiment analysis:

1) Slang words translation: The most common Tweets often contain a lot of slang words (e.g. lol, omg). These words are usually very important for sentiment analysis, but may not be included in root sentiment lexicons. Since the sentiment analysis is based on sentiment lexicon, therefore we are converting these all slang words into their standard forms using the Internet Slang Word Dictionary and then re-add them to the tweets.

2) Non-English tweets filtering: Since the sentiment analysis tools to be used only work for

English texts, we remove all non-English tweets in advance as these non English words doesn't have meaning for sentiment. A tweet could be treated as non-English tweet if more than 20 percent of its words (after slang words translation) do not appear in the GNU Aspell English Dictionary.

3) URL removal: A lots of users may include various URLs in their tweets. These URLs may complicate our sentiment analysis process. So we decide to remove URLs from tweets. (100%)

B. Sentiment Label Assignment For assigning sentiment labels for each tweet more confidently, we sort lexicons again to two state-of-the-art sentiment analysis tools. One is the SentiStrength3 tool [8]. This tool is based on the LIWC [10] sentiment lexicon. It works in the following way: first assign a sentiment score to each word in the text according to the sentiment lexicon; then choose the maximum positive score and the maximum negative score among those of all individual words in the text; compute the sum of the maximum positive sentimental score and the maximum negative sentimental score, denoted as Final sentimental Score; finally, use the sign of Final Score to indicate whether a tweet is positive, negative or it is neutral.

Lexicon based Techniques

In unsupervised technique, classification is done by comparing the features of a given text against sentiment lexicons whose sentiment values are determined prior to their use. Sentiment lexicon contains lists of words and expressions used to express people's subjective feelings and opinions. For example, start with positive and negative word lexicons, analyse the document for which sentiment need to find. Then if the document has more positive word lexicons, it is positive, otherwise it is negative. The lexicon based techniques to Sentiment analysis is unsupervised learning because it does not require prior training in order to classify the data.

The steps of the lexicon based techniques are below

1. Preprocess each raw tweet text (i.e. remove HTML tags, noisy characters).

2. Initialize the total text sentiment score: $s = 0$.

3. Tokenize text. For each token, check if tokens are present in a sentiment dictionary of training set.

- (a) If token is present in dictionary,

- i. If token is positive, then $s = s + w$.

- ii. If token is negative, then $s = s - w$.

4. Look at aggregate text sentiment score s ,

- (a) If $s > \text{threshold}$, then classify the text as positive

- (b) If $s < \text{threshold}$, then classify the text as negative.

V. CONCLUSIONS

Overall, we conclude that social network based behavioral analysis parameters can increase the prediction accuracy. However, presence of all the entities in unbiased and equal manner is necessary to provide accurate results. In this paper, we investigated the problem of analyzing public sentiment variations and finding the possible reasons causing these variations. we proposed two Latent Dirichlet Allocation (LDA) based models, Foreground and Background LDA (FB-LDA) and Reason Candidate and Background LDA (RCB-LDA). These foreground topics can give potential interpretations of the sentiment variations. we have selected the descriptive tweets for foreground topics and develop another generative model called Reason Candidate and Background LDA (RCB-LDA) to rank them with respect to their —popularity— within the variation period. The FB-LDA model can filter out background topics and then extract foreground topics to reveal possible reasons. To give a more spontaneous, representation of the RCB-LDA model which can rank a set of reason candidates expressed in natural language to provide sentence-level reasons. The proposed models are general: they can be used to discover special topics or aspects in one text collection in comparison with another

background text collection. Also our proposed models evaluated on real Twitter data.

REFERENCES

- [1] B. O'Connor, R. Balasubramanyan, B. R. Routledge, and N. A. Smith, "From tweets to polls: Linking text sentiment to public opinion time series," in Proc. 4th Int. AAAI Conf. Weblogs SocialMedia, Washington, DC, USA, 2010.
- [2] Bo Pang. Lilliam Lee, "Seeing Stars: Exploiting class relationships for sentiment categorization with respect to rating scales", 2002
- [3] J. Bollen, H. Mao, and A. Pepe, "Modeling public mood and emotion: Twitter sentiment and socio-economic phenomena" ,in Proc.5th Int. AAAI Conf. Weblogs Social Media, Barcelona, Spain, 2011.
- [4] J. Bollen, H. Mao, and X. Zeng, —Twitter mood predicts the stock market," J. Comput. Sci., vol. 2, no. 1, pp. 1–8, Mar. 2011.
- [5] Bernard J. Jansen, Mimi Zhang, Kate Sobel and AbdurChowdury,□Micro-blogging as online word of mouth branding", 27th International Conference Extended Abstracts on Human Factors in Computing Systems, New York, 2009, pages 3859-3862.
- [6] J.C. Na, H. Sui, C. Khoo, S. Chan, and Y. Zhou.□Effectiveness of simple linguistic processing in automatic sentiment classification of product reviews",Advances In Knowledge and organization, 2004, pages 49-54.
- [7] T. Minka and J. Lafferty, —Expectation-propagation for the generative aspect model□, in Proc. 18th Conf. UAI, San Francisco, CA, USA, 2002.
- [8] M. Thelwall, K. Buckley, G. Paltoglou, D. Cai, and A. Kappas, —Sentiment strength detection in short informal text□ „J. Amer. Soc.Inform. Sci. Technol., vol. 61, no. 12, pp. 2544–2558, 2010.
- [9] A. Go, R. Bhayani, and L. Huang, —Twitter sentiment classification using distant supervision□, CS224N Project Rep., Stanford: 1–12, 2009.
- [10] Y. Tausczik and J. Pennebaker, —The psychological meaning of words: Liwc and computerized text analysis methods□ , J. Lang.Soc. Psychol., vol. 29, no. 1, pp. 24–54, 2010.