# Survey of Public Sentiment Interpretation on Twitter

Ms. Devaki V. Ingule[1], Prof.Gyankamal J. Chhajed[2]

[1]Student, M. E. Computer Engineering, VPCOE, Baramati, SavitribaiPhulePuneUniversity, India

*devaki.ingule@gmail.com*

[2]Assistant Professor, Computer Engineering Department, VPCOE, Baramati, Savitribai Phule PuneUniversity, India

*gjchhajed@gmail.com*

**Abstract**—Opinion mining and sentiment analysis is a Natural Language Processing task. Number of Users shared what they think on microblog services. Twitter is important platform for follow sentiment of public which is a very challenging problem. Public sentiment analysis is a very essentialto explore, analyze and organize user's views for better decision making. Sentiment analysis is process of identifying positive and negative opinions, emotions and evaluations in text. This is useful for consumers for research the sentiment of products before they actually purchase, or companies that want to monitor the public sentiment of their brands. In this paper we have reviewed and analyzed number of techniques for public sentiments analysis and their classification.

**Keywords**—Twitter,Public sentiment, Emerging Topic Mining, Event tracking, sentiment classification, supervisedmachine learningmethods, Correlation between Tweets and Events.

## INTRODUCTION

Mining public sentiments and analysis of them on twitter data has provided easy way to expose public opinion, which helps for decision making in various domains. Twitter is important and popular platforms for peoples interaction. By using twitter platform number of users share their views and opinions. For making important decision it is necessary to mine public opinions and to find reasons behind variation of sentiments is valuable. For example, acompany can analyze opinions of public for obtainingusers feedback about its products in tweets. In general, opinion mining helps to collect information about the positive and negative aspects of a particular topic. Finally, the positive and highly scored opinions obtained about a particular product are recommended to the user. In this paper number of different methods are used for analysis of public sentiments, opinion is classified in various approaches on text using some Supervised machine learning algorithm like Maximum Entropy classification, Support Vector Machines, Naive Bayes. Combination of two state-of-the-art sentiment analysis tools is used for obtaining sentiment information.

## 1. SENTIMENT ANALYSIS

Sentiment Analysis, analyze the opinions which are extracted from various sources like the comments on forums, reviews about products, different policies or topic related with social networking sites and tweets.

O'Connor et al.studied on analyzing sentiments of public share on Twitter. This is the first work in microblogging services to interpret the variations in sentiment.

**A]** Pang et al. [1] studied the existing methods on analysis of sentiments in details i.e.supervised machine learning methods.

**Advantages:**

1. Machine learning methods minimize the structural risks.
2. For predicts sentiment of documents Supervised machine learning approach are used.

**Disadvantages:**

1. It cannot analyze possible reasons behind public sentiments.
2. Supervised methods demand large amounts of labeled training data that are very expensive.
3. It may fail when training data are insufficient.

**B]** M.Hu and B. Liu [2] also works on mining and summarizing customer reviews based on data mining and natural language processing methods. This paper develops different novel techniques to summarize reviews of customers. It summarizes reviews by following three ways:

1. It mine features of products on which customers have been commented.
2. It Identifies opinions and decides whether it is positive or negative in each review.
3. It also summarizing the results of opinions.

**Advantages:**

1. It predicts movie sales and elections so it's easy to make decisions.
2. It provides a feature-based summary of a large number of customer reviews of a product sold online.
3. Summarizing the reviews is useful to common shoppers, also to product manufacturers.

**Disadvantages:**

1. It cannot determine the opinions strength.
2. It also not expressed opinions with adverbs, verbs and nouns.

**C]** W.zhang et al. [3] conducted detailed study of opinions retrieval from blogs. In this paper, they have presented a three-component opinion retrieval algorithm. The first component is an information retrieval module. The second one classifies documents into opinionative and no opinionative documents, and keeps the former. The third component ensures that the opinions are related to the query, and ranks the documents in certain order.

Meng et al.collected opinions in Twitter for entities by mining hash-tags to conclude the presence of entities and sentiments from each tweets.

**Advantages:**

1. It mines hash tags for opinions.
2. It has higher performance than a state-of-art opinion retrieval system.

**Disadvantages:**

1. It cannot handle more general writings and crossing domains.
2. It cannot select detail features.
3. It cannot implement better NEAR operator.
4. Because of all tweets cannot contain hash tags, it is difficult to gain sufficient coverage for an events.

Sentiment classification is one of the best applications of sentiment analysis which classifies text into number of categories.

**D]** Li.Jiang et.al [4] described the target-dependent Twitter sentiment Classification. The state-of-the-art technique solves the problem of target dependent classification which adopts the target-independent strategy. To given target it always assign irrelevant sentiments .For agiven query, it classifies the sentiments of the tweets according to there categories whether they contain positive, or negative or any neutral sentiments about the query. Here the query considers as the target of the sentiments.

**Advantages:**

1. This technique gives the high performance for target-dependent twitter sentiment classification.
2. By using words which are connected to the given target it incorporates syntactic features that are generated using words syntactically to decide sentiment is about the given target or not.

**Disadvantages:**

1. It is not search out the relations between a target and its extended targets.
2. It is not able to explore relations between twitter accounts that classifies the sentiments.
3. It decreases the performance of very short and ambiguous tweets.

## 2. EVENT DETECTION AND TRACKING

Events are the good reasons behind the variations of sentiments related to target. By tracking the events, and detecting sentiment analysis about events and also tracking variations in sentiments and finding reasons for changes the sentiments completes this task.

**A]** Leskovec et al.[5] proposed work on tracking memes, for example quoted phrases and sentences. This work offers some of analysis of the global news cycle and the dynamics of information propagation between mainstream and social media. It identifies short distinctive phrases that travel relatively intact through online text as it evolves over time.

**Advantages:**

1. It provides temporal relationships such as the possibility of employing a type of two-species predator-prey model with blogs and the news media as the two interacting participants.

**Disadvantages:**

1. It is applicable only for representative events, such as biggest event in whole twitter message stream.
2. Fine grained events can be detected very hardly.

**B]** B-S Lee and J.Weng [6] concentrate on detection of events through analysis of the contents which are published in Twitter. This paper proposes EDCoW that is Event Detection with Clustering of Wavelet based on Signals for detecting events.

**Advantages:**

1. EDCoW (Event Detection with Clustering of Wavelet-based Signals) signal independently treats each word.
2. EDCoW (Event Detection with Clustering of Wavelet-based Signals) achieves good performance.
3. This could possibly contribute to study the temporal evolution of event.

**Disadvantages:**

1. It cannot contribute relationship among users to event detection.
2. This design of EDCoW not applicable on time lag when it computing for the cross correlation between a pair of words.

## 3. DATA VISUALIZATION

**A]** D.Tao et.al [7] deeply studied subspace learning algorithms and ranking. Retrieving of images from large databases is very active research field today. For retrieving images content based image retrieval (CBIR) technique is used. It is related by semantically to query of user from an collected database of images.SVM classifier always unstable for a smaller size training set. SVMRF becomes poor if there are number of samples of positive feedback are small.SVM has also an problem of over fitting.

**Advantages:**

1. It improves the performance of relevance feedback.
2. SVM classifier always unstable on a smaller size of training set to address this it develops an asymmetric bagging-based SVM (AB-SVM).
3. For over fitting problem it used the random subspace method and SVM together for relevance feedback.

**Disadvantages:**

1. It cannot use tested tuning method and not select the parameters of kernel based algorithms.
2. These works are not useful for text data especially for noisy text data.
3. Because of explicit queries are not present in task, ranking methods cannot solve the reason mining task.

## 4. CORRELATION BETWEEN TWEETS AND EVENTS

**A]** T.Sakaki et al.[8] developed novel models to map tweets in a public segmentation. They detect real-time events in Twitter such as earthquakes. Theyalso proposes an algorithm for monitoring tweets detect event. Each Twitter user is considering as a sensor.Kalman filtering and particle filtering are used for estimation of location.

**Advantages:**

1. Earthquakes are detected by this system and sends e-mails to users who are registered users.
2. Kalman filtering and particle filtering detects and provides estimation for location.

**Disadvantages:**

1. It cannot detect multiple event occurrences.
2. It cannot provide advanced algorithms for query expansion.
3. There was only one target event at a time.


**B]** Y.Hu et.al [9] proposed a joint statistical model ET-LDA that characterizes topical effects between an event and its related Twitter feeds. This model enables the topic modeling of the event and the segmentation of the each event or tweet.

**Advantages:**

1. It extracts a variety of dimensions such as sentiment and polarity.
2. It describes the temporal correlation between overall tweets.

**Disadvantages:**

1. They model each tweet as a multinomial mixture of all events, which is obviously unreasonable due to short lengths of tweets.


**C]** Chakrabarti and Punera [10] have described a variant of Hidden Markov Models in event summarization from tweets. It gets an intermediate representation for a sequence of tweets relevant for an event. In this paper use of sophisticated techniques to summarize the relevant tweets are used for some highly structured and recurring events. Hidden Markov Models gives the hidden events.

**Advantages:**

1. It provides benefits for existing query matching technologies.
2. It works well for one-shot events such as earthquakes.
3. It tackled the problem of constructing real time summaries of events.
4. It learns an underlying hidden state representation of an event.

**Disadvantages:**

1. It does not use the continuous time stamps present in tweets.
2. It is not possible to gets minimal set of tweets which are relevant to an event.
3. It cannot provide a summary of long occurrences and unpredictable events.
4. In above novel model noises and background topics cannot be eliminated.


**[D]**Shulong Tan et.al [11] proposes LDA model for interpreting public sentiment variations on Twitter. WhereLatent Dirichlet Allocation (LDA)  propose two models 1.Foreground and Background LDA (FB-LDA) and 2.Reason Candidate and Background LDA (RCB-LDA) in which FB-LDA model can remove background topics and then extract foreground topics to show possible reasons and Reason Candidate and Background LDA(RCB-LDA) model provides ranking them with respect to their popularity within the variation period. RCB-LDA model also finding the correlation between tweets and their events.

## CONCLUSION

This survey discusses various approaches to Opinion Mining and Sentiment Classification. It provides a detailed view of different applications and potential challenges of Sentiment Classification. Some of the machine learning techniques like Naive Bayes, Maximum Entropy and Support Vector Machines and their pros and cons has been discussed. The emerging topics are related to the actual or genuine reasons behind the variations are very important. So emerging topics will consider as possible reasons behind variations.It is necessary to interpret sentiment variation and finding the reasons behind them for overcome above limitations.One of the technique that is Latent Dirichlet Allocation (LDA) which propose two models such as Foreground and Background LDA (FB-LDA) and Reason Candidate and Background LDA (RCB-LDA)willinterprets sentiment variations.It concludes that these models can mine number of reasons behind variations of public sentiments.

## REFERENCES:

[1] B. Pang and L. Lee, "Opinion mining and sentiment analysis,"*Found. Trends Inform. Retrieval*, vol. 2, no. (1–2), pp. 1–135, 2008.

[2] M. Hu and B. Liu, "Mining and summarizing customer reviews,"in *Proc. 10th ACM SIGKDD*, Washington, DC, USA, 2004.

[3] W. Zhang, C. Yu, and W. Meng, "Opinion retrieval from blogs,"in *Proc. 16th ACM CIKM*, Lisbon, Portugal, 2007.

[4] L. Jiang, M. Yu, M. Zhou, X. Liu, and T. Zhao, "Target-dependent twitter sentiment classification," in *Proc. 49th HLT*, Portland, OR,USA, 2011.

[5] J. Leskovec, L. Backstrom, and J. Kleinberg, "Meme-tracking and the dynamics of the news cycle," in *Proc. 15th ACM SIGKDD*,Paris, France, 2009.

*[6]* J. Weng and B.-S. Lee, "Event detection in twitter," in *Proc. 5thInt. AAAI Conf. Weblogs Social Media*, Barcelona, Spain, 2011

[7] D. Tao, X. Tang, X. Li, and X. Wu, "Asymmetric bagging and random subspace for support vector machines-based relevance feedback in image retrieval," *IEEE Trans. Patt. Anal. Mach. Intell.*,vol. 28, no. 7, pp. 1088–1099, Jul. 2006

[8] T. Sakaki, M. Okazaki, and Y. Matsuo, "Earthquake shakes twitter users: Real-time event detection by social sensors," in *Proc. 19th Int. Conf. WWW*, Raleigh, NC, USA, 2010.

[9] Y. Hu, A. John, F. Wang, and D. D. Seligmann, "Et-lda: Joint topic modeling for aligning events andtheir twitter feedback," in *Proc.26th AAAI Conf. Artif. Intell.*, Vancouver, BC, Canada, 2012.

[10] D. Chakrabarti and K. Punera, "Event summarization using tweets," in *Proc. 5th Int. AAAI Conf. Weblogs Social Media*, Barcelona, Spain, 2011.

[11] Shulong Tan, Yang Li, Huan Sun, Ziyu Guan, Xifeng Yan, "Interpreting the Public Sentiment Variations on Twitter," IEEE Transactions on Knowledge and Data Engineering, VOL. 26, NO.5, MAY 2014.

[12] L. Zhuang, F. Jing, X. Zhu, and L. Zhang, "Movie review mining and summarization," in *Proc. 15th ACM Int. Conf. Inform.Knowl.Manage.*, Arlington, TX, USA, 2006.

[13] X.Wang, F.Wei, X. Liu, M. Zhou, and M. Zhang, "Topic sentiment analysis in twitter: A graph-based hashtag sentiment classification approach," in *Proc. 20th ACM CIKM*, Glasgow, Scotland, 2011.

[14] 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