

Clustering Based real time Sentiment Analysis of Big Data

^{#1}Prof. R. B. Pawar, ^{#2}Yogesh Jadhav, ^{#3}Tejaswini Jagadale, ^{#4}Shraddha Kale,
^{#5}Harsha Valwani



¹jadhavyogesh611@gmail.com,
²tejaswinijagadale1@gmail.com,
³shraddhakale119@gmail.com,
⁴harshavalwani1998@gmail.com

^{#12345}Department Of Information Technology, AVCOE Sangamner,
Savitribai Phule Pune University, Maharashtra.

ABSTRACT

The massive growth of social media has given fresh impetus to the field of sentiment analysis. It is the computational treatment of sentiments in a text. Through this paper the authors have given an overview of the different sentiment analysis techniques and tools currently being employed to gauge the sentiments of social media data. The real challenge with real-time stream data processing is that it is impossible to store instances of data, and therefore online analytical algorithms are utilized. To perform real-time analytics, pre-processing of data should be performed in a way that only a short summary of stream is stored in main memory. The application of opinion mining and sentiment analysis (OMSA) in the era of big data have been used a useful way in categorize the opinion into different sentiment and in general evaluating the mood of the public.

Keyword: Opinion mining, sentiment analysis, big data, online social network.

ARTICLE INFO

Article History

Received: 20th March 2019

Received in revised form :

20th March 2019

Accepted: 22nd March 2019

Published online :

23rd March 2019

I. INTRODUCTION

Sentimental analysis on big data is useful to analysis of the product feedback. Now a day's people are posting on social media about the product like twitter. So it is useful to get analysis on online data. In sentimental analysis done by static on small data has issue like the feedback is limited data. However, big data is generally defined through the key characteristics of volume, variety, and velocity [8].

1. Volume represents the quantity of data that utilizes massive storage space or entails sizeable and considerable number of records [9]. For instance, WalMart manages to store 2.5 petabytes of information, Tesco generated 1.5 billion new items of data on a monthly basis, and Dell developed a database which can handle 1.5 million sales and advertisement records [3], [10].

2. Variety refers to the data generated from a range of sources and in varying formats [9]. The sources can be in the form of sensors, social media sites, web technologies, mobile phones, etc. The data format can include web logs, unstructured data like audio, videos, images and sensory data from RFID devices, or other smart sensors.

3. Velocity indicates the frequency with which data is produced from different sources [9]. The data can be generated occasionally, frequently, and/or on a real-time basis.

Current solutions and studies in data stream sentiment analysis are limited to perform sentiment analysis in an off-line approach on a sample of stored stream data. While this approach can work in some cases, it is not applicable in the real-time case. In addition, real-time sentiment analysis tools such as MOA and Rapid Miner exist, however they are uniprocessor solutions and they cannot be scaled for an efficient usage in a network nor a cluster. Since in big data scenarios, the volume of data rises drastically after some period of analysis, this causes uniprocessor solutions to perform slower over time. As a result, processing time per instance of data becomes higher and instances get lost in a stream. This affects the learning curve and accuracy measures due to less available data for training and can introduce high costs to such solutions. Sentinel relies on distributed architecture and distributed learner's to solve this shortcoming of available solutions for real-time sentiment analysis in social media.

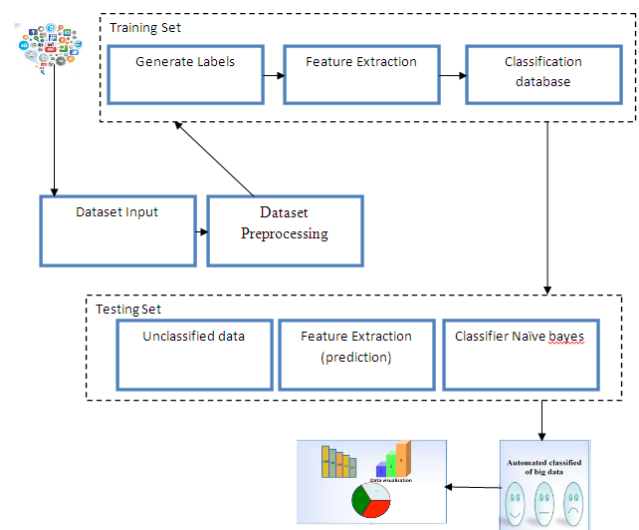
II. LITERATURE SURVEY

An opinion refers to a person's or group's sentiment or views, emotions, and attitudes about a product, service, occasion, or other topic present in the environment. Like sentiment analysis, opinion mining is also grounded on the algorithmic technique [21]. Covering a huge variety of public opinions, [22] have argued that opinion can be classified into three main types: regular opinions, which refer to a single entity only; comparative opinions, which compare or contrast more than one entity; and suggestive opinions, which suggest a single or multiple entities. The regular opinion is mainly used to identify a positive or negative outlook of a particular product [23]. On the other hand, comparative opinions help in elucidating the association among multiple entities and are mainly used for competitive intelligence [23]. However, there is a dearth of literature concerning the identification of comparative sentences that is being used for the comparison of multiple entities. Recently, suggestive review has been introduced in the field of opinion mining [24]. The extraction of these types of opinions from text can be utilized for various application areas in the field of business, engineering, medical science, and e-learning. It can be offered to various online communities for their assistance as well [22]. Similarly, private statements of individuals are called sentiments, which comprise thoughts, opinions, attitudes, views, judgments, and feelings. These are commonly gathered by conventional scientific methods [19, 20]. [25] pronounced the feelings that are expressed in language by using subjective expression. The sentiments can be analysed through the machine learning technique, which can be further classified into supervised and unsupervised, using a lexicon based approach, using a keyword, and using a concept-based technique [26]. Recently, research on sentiment analysis has focused on multiple modalities such as in speech and video as opposed to earlier work that focused on unimodality related to text [27, 28]. Sentiment analysis tackles many NLP subtasks, including aspect extraction [29], subjectivity detection [30], named entity recognition, and sarcasm detection [31]. In most cases, the main objective of sentiment analysis is to unearth people's opinions to gain meaningful insight about products or services. Its aim is to exhibit useful information to both customers and manufacturers.

It is established that both manufacturer and customers look upon summarized opinions instead of detailed reviews. Hence the opinions that are categorized on positive, negative, or neutral sentiments are useful for both parties in making the right call [32]. Despite the large number of studies on opinion mining and sentiment analysis techniques, the impact they have on people has been less explored. There has been great emphasis on the techniques used and less on how people can benefit from the findings. Hence, this study aims to investigate the human element in opinion mining and sentiment analysis. research. To achieve this aim, we will systematically review the relevant literatures that have employed both approaches. The study offers several contributions. The first

and foremost significance of this study is to refocus the study of opinion mining and sentiment analysis to both technical and nontechnical challenges. Secondly, it places emphasis on the areas of opportunity by looking at the trends of application coverage that would offer some potential areas for research. Thirdly, the paper presents information on different datasets that were used in opinion mining and sentiment analysis studies that future researchers could use in their research. The remaining section of the paper is organized as follows. First, the study talks about the method employed to achieve the research objectives. Then we present the findings of the study. Next, the paper highlights the commonly used dataset in literature.

III. METHOD DESCRIPTION



Process:

Preprocessing Of Data

Feature Extraction

Probability Calculation

Clustering data

Classification

Preprocessing

To pre-process the data in order to remove noise and unrelated content. It should be followed by the construction and evaluation of the sentiment analysis model (based on keywords, lexicon, or machine learning methods).

Feature Extraction

Feature extraction is an important part of building an effective machine learning method in which the textual posts are transformed into valuable word features by using various feature engineering approaches. Feature extracting is one the most important steps of constructing effective classifiers. The accomplishment or failure of the sentiment classification model is intensely dependent on the features quality. If the extracted features relate well with the

sentiment polarity and can provide discrimination power between positive and negative, then classification will be more precise.

Probability Calculation

System is then used to calculate the probabilities of each class. The conditional probabilities of each class are calculated using single instance from the test set. Posterior probabilities of each class are then calculated. This process is undergone on each instance of the dataset.

$$p(w_n | c) = \frac{m_{w_n c}}{\sum_{n=1}^N m_{w_n c}}$$

Clustering Data

Sentiment clustering is composed of two steps: grouping data based on timestamp and clustering sentiment based on topic. First step aims at clustering the entire dataset into different date categories with day and week labels using personalized dataset temporal clustering algorithm. Second step clusters sentiment using kmeans clustering algorithm with bag-of-words term weighting method and distance matrix in accordance with topic.

Classification

One popular way to implement multi-label classifier is to transform the multi-label classification problem into multiple single-label classification problems. One simple transformation method is called one-versus-all or binary relevance. The basic concept is to assume independence among categories, and train a binary classifier for each category. All kinds of binary classifier can be transformed to multi-label classifier using the one-versus-all heuristic. The basic procedures of the multi-label Naïve Bayes classifier.

Naïve Bayes Pseudo Code

1. Let's call vs the list of vertices in the graph. $len(vs)$ is the length. $vs[i]$ is the i th vertex.
2. Let's assume we have a univariate and binary scenario, i.e., $vs[i].class$ is either 0 or 1 and there is no other given feature of a node.
3. Let's assume we run a local classifier before so that every node has an initial label, which are calculated by the local classifier. I am only interested in with the relational classifier part.
4. Let's call v the vertex we are trying to predict, and $v.neighbors()$ is the list of vertices which are neighbors of v .
5. Let's assume all the edge weights are 1.

IV. CONCLUSION

The main contribution of this paper is the design and development of a real time big data stream analytic framework; providing a foundation for an infrastructure of real time sentiment analysis on big text streams. Our

framework is proven to be an efficient, scalable tool to extract, score and analyze opinions on user generated text streams per user given topics in real time or near real time.

V. ACKNOWLEDGMENT

In this paper, we thank our guide Prof. R.B.Pawar for the valuable information, feedback and comments and have helped us immensely and contributed their time and idea to this research.

REFERENCES

- [1] R. Addo-tenkorang and P. T. Helo, "Computers & Industrial Engineering Big data applications in operations / supply-chain management : A literature review," *Comput. Ind. Eng.*, vol. 101, pp. 528–543, 2016.
- [2] A. Zaslavsky, C. Perera, and D. Georgakopoulos, "Sensing as a Service and Big Data," 2015.
- [3] J. Manyika et al., "Big data: The next frontier for innovation, competition, and productivity," 2011.
- [4] H. Sundmaeker, P. Guillemin, P. Friess, and S. Woelfflé, "Vision and challenges for realising the Internet of Things," *Clust. Eur. Res. Proj. Internet Things*, Eur. Commission, vol. 3, no. 3, pp. 34–36, 2010.
- [5] A. N. E. W. Approach and T. O. Analysis, "Unstructured data : A big deal in big data Deep dive : Analytics."
- [6] R. L. Villars, C. W. Olofson and M. Eastwood, "Big Data: What It Is and Why You Should Care Information," White Paper: IDC, June 2011 .
- [7] M. Khoso, "How much data is produced every day," 2016.
- [8] P. Zikopoulos, C. Eaton, and others, *Understanding big data: Analytics for enterprise class hadoop and streaming data*. McGraw-Hill Osborne Media, 2011.
- [9] P. Russom and p-oIStHSeNrs:23"9B4i-g82d8a0ta analytics," TDWI best Pract. report, fourth Quart., vol. 19, p. 40, 2011.
- [10] T. H. Davenport, "Big Data in Big Companies," *Int. Inst. Analytics*, 3. May, 2013.
- [11] A. Global and O. Megatrends, "Big Data Analytics in Supply Chain : Hype or Here to Stay ? Big data analytics," pp. 1–20.
- [12] F. Provost and T. Fawcett, "Data Science and its relationship data-driven decision making," *Big Data*, vol. 1, no. 1, pp. 51–59, 2013.
- [13] H. Chen and V. C. Storey, "Business Intelligence and Analytics: From Big Data to Big Impact." *MIS Quarterly*, vol. 36, no. 4, pp. 1165– 1188, 2012..
- [14] A. Misra, A. Sharma, P. Gulia, and A. Bana, "Big Data : Challenges and Opportunities," no. 2, pp. 41–42, 2014.
- [15] M. Batty et al., "Smart cities of the future," *Eur. Phys. J. Spec. Top.*, vol. 214, no. 1, pp. 481– 518, 2012.

- [16] R. E. Bryant, R. H. Katz, and E. D. Lazowska, "Big-Data Computing : Creating revolutionary breakthroughs in commerce , science , and society Motivation : Our Data- Driven World," 2008.
- [17] F. R. Lucini et al., "Text mining approach to predict hospital admissions using early medical records from the emergency department," *Int. J. Med. Inform.*, vol. 100, pp. 1–8, 2017.
- [18] Z. Khan, Z. Khan, T. Vorley, and T. Vorley, "Big data text analytics: an enabler of knowledge management," *J. Knowl. Manag.*, vol. 21, no. 1, pp. 18–34, 2017.
- [19] T. T. Thet, J.-C. Na, and C. S. G. Khoo, "Aspectbased sentiment analysis of movie reviews on discussion boards," *J. Inf. Sci.*, vol. 36, no. 6, pp. 823–848, 2010.
- [20] H. Yu and V. Hatzivassiloglou, "Towards answering opinion questions: Separating facts from opinions and identifying the polarity of opinion sentences," in *Proceedings of the 2003 conference on Empirical methods in natural language processing*, 2003, pp. 129–136.
- [21] R. Piryani, D. Madhavi, and V. K. Singh, "Analytical mapping of opinion mining and sentiment analysis research during 2000--2015," *Inf. Process. Manag.*, vol. 53, no. 1, pp. 122–150, 2017.
- [22] A. Qazi, A. Tamjidyamcholo, R. G. Raj, G. Hardaker, and C. Standing, "Assessing consumers' satisfaction and expectations through online opinions: Expectation and disconfirmation approach," *Comput. Human Behav.*, vol. 75, pp. 450–460, 2017.
- [23] N. Jindal and B. Liu, "Identifying comparative sentences in text documents," in *Proceedings of the 29th annual international ACM SIGIR conference on Research and development in information retrieval*, 2006, pp. 244–251.
- [24] A. Qazi, R. G. Raj, M. Tahir, E. Cambria, and K. B. S. Syed, "Enhancing business intelligence by means of suggestive reviews," *Sci. World J.*, vol. 2014, 2014.
- [25] M. Quigley, *Encyclopedia of information ethics and security*. IGI Global, 2007.
- [26] E. Cambria, "Affective computing and sentiment analysis," *IEEE Intell. Syst.*, vol. 31, no. 2, pp. 102–107, 2016.
- [27] S. Poria, E. Cambria, N. Howard, G.-B. Huang, and A. Hussain, "Fusing audio, visual and textual clues for sentiment analysis from multimodal content," *Neurocomputing*, vol. 174, pp. 50–59, 2016.
- [28] S. Poria, I. Chaturvedi, E. Cambria, and A. Hussain, "Convolutional MKL based multimodal emotion recognition and sentiment analysis," in *Data Mining (ICDM), 2016 IEEE 16th International Conference on*, 2016, pp. 439–448.
- [29] S. Poria, E. Cambria, and A. Gelbukh, "Aspect extraction for opinion mining with a deep convolutional neural network," *Knowledge-Based Syst.*, vol. 108, pp. 42–49, 2016.
- [30] I. Chaturvedi, E. Cambria, S. Poria, and R. Bajpai, "Bayesian Deep Convolution Belief Networks for Subjectivity Detection," in *Data Mining Workshops (ICDMW), 2016 IEEE 16th International Conference on*, 2016, pp. 916–923.