

Application of Machine Learning Techniques to Sentiment Analysis

Anuja P Jain

Computer Science and Engineering, Mtech
Gogte Institute of Technology
Belgaum, India
jainanu04@gmail.com

Asst. Prof Padma Dandannavar

Computer Science and Engineering
Gogte Institute of Technology
Belgaum, India
padmad@git.edu

Abstract—Today, we live in a ‘data age’. Due to rapid increase in the amount of user-generated data on social media platforms like Twitter, several opportunities and new open doors have been prompted for organizations that endeavour hard to keep a track on customer reviews and opinions about their products. Twitter is a huge fast emergent micro-blogging social networking platform for users to express their views about politics, products sports etc. These views are useful for businesses, government and individuals. Hence, tweets can be used as a valuable source for mining public’s opinion. Sentiment analysis is a process of automatically identifying whether a user-generated text expresses positive, negative or neutral opinion about an entity (i.e. product, people, topic, event etc). The objective of this paper is to give step-by-step detail about the process of sentiment analysis on twitter data using machine learning. This paper also provides details of proposed approach for sentiment analysis. This work proposes a Text analysis framework for twitter data using Apache spark and hence is more flexible, fast and scalable. Naïve Bayes and Decision trees machine learning algorithms are used for sentiment analysis in the proposed framework.

Keywords—sentiment analysis; machine learning; Natural Language Processing; twitter;

I. INTRODUCTION

Twitter is a tremendous, quick, rising, prevalent, smaller scale blogging interpersonal interaction platform for users to express their perspectives about governmental issues, product items, sports and so forth. Here, users send messages (a.k.a., tweets) to a network of contacts from a wide variety of devices. A tweet is a text-based post and only has 140 characters, which is approximately the length of a typical newspaper headline and subhead [4]. Twitter is a “what’s-happening-right-now” social network and hence tweets are valuable sources for businesses, government and individuals to determine public’s opinion or sentiment about an entity (product, people, topic, event etc). But, the volume of tweets produced by Twitter everyday is very vast (21 million tweets per hour, as measured in 2015). Hence there is a need to automate the process of sentiment analysis so as to ease the tasks of determining public’s opinions without having to read millions of tweets manually. This process of analyzing and summarizing the opinions expressed in these huge opinionated user generated data is usually called Sentiment Analysis or Opinion Mining which is a very interesting and popular domain for researchers nowadays.

Sentiment analysis is a process of automatically identifying whether a user-generated text expresses positive, negative or neutral opinion about an entity (i.e. product, people, topic, event etc). Sentiment classification can be done at Document level, Sentence level and Aspect or Feature level [1]. In Document level the whole document is used as a basic information unit to classify it either into positive or negative class. Sentence level sentiment classification classifies each sentence first as subjective or objective and then classifies into positive, negative or neutral class. There is no much difference between the above two methods as sentence is just a short document. Aspect or Feature level sentiment classification deals with identifying and extracting product features from the source data [1].

There are several approaches for sentiment analysis: *Machine learning based approach* (ML) uses several machine learning algorithms (supervised or unsupervised algorithms) to classify data. *Lexicon based approach* uses a dictionary containing positive and negative words to determine the sentiment polarity. *Hybrid based approach* uses a combination of both ML and lexicon based approach for classification.

The organization of paper is as follows: Section 1 gives an overview about various sentiment classification techniques. Section 2 explains the various steps needed for sentiment analysis using machine learning and NLP. Section 3 provides an outline of machine learning algorithms for sentiment classification.

II. SENTIMENT CLASSIFICATION TECHNIQUES

A. Lexicon based approach

It is based on finding the opinion lexicon for calculating the sentiment for a given text. It deals with counting the number of positive and negative words in the text. If the text consists of more positive words, the text is assigned a positive score. If there are more number of negative words the text is assigned a negative score. If the text contains equal number of positive and negative words then it is assigned a neutral score. To determine whether a word is positive or negative an opinion lexicon (positive and negative opinion words) is built. There are several approaches to compile and build an opinion lexicon [2]:

- Dictionary based approach: A small set of opinion words is collected manually with known orientations [2]. Then, synonyms and antonyms of these words are searched in corpora like WordNet or thesaurus and added to the set. The set gradually grows until no new words are found. This approach has a disadvantage that the strength of the sentiment classification depends on the size of the dictionary. As the size of the dictionary grows this approach becomes more erroneous.
- Corpus based approach: They depend on large corpora for syntactic and semantic patterns of opinion words. The words that are generated are context specific and may require a huge labelled dataset.

B. Machine Learning based approach

Here, two sets of documents are needed: training and a test set. A supervised learning classifier uses the training set to learn and train itself with respect to the differentiating attributes of text, and the performance of the classifier is tested using test dataset. Several machine learning algorithms like Maximum Entropy (ME), Naive Baye's (NB) and Support Vector Machines (SVM) are usually used for classification of text (tweets).

Machine Learning for sentiment analysis starts with collection of dataset containing labelled tweets. This dataset might be boisterous and subsequently should be pre handled utilizing various Natural Language processing (NLP) techniques. Then features that are relevant for sentiment analysis need to be extracted and finally the classifier is trained and tested on unseen data. These algorithms are explained in detail in section 4.

The various steps performed for sentiment analysis using machine learning is shown in figure 1. Machine Learning for sentiment analysis starts with collection of dataset containing labelled tweets. This data may be noisy and hence needs to be pre-processed using various Natural Language processing (NLP) techniques. Then features that are relevant for sentiment analysis need to be extracted and finally the classifier is trained and tested on unseen data. These steps are explained in detail in section 3.

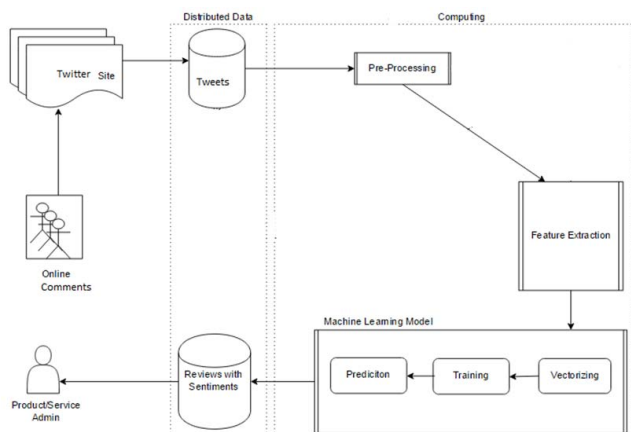


Fig 1: Workflow for twitter sentiment analysis using machine learning algorithms.

C. Hybrid approach

In order to improve sentiment classification performance few research techniques suggest using a combination of both lexicon based and machine learning techniques. The main advantage of this hybrid approach is that we can attain best of both world. The lexicon/learning combination has proven to improve accuracy. As mentioned in [5], lexicon based approach have high precision and low recall. Hence combining it with a machine learning classifier can improve the recall and accuracy of the algorithm.

This paper focuses on sentiment analysis using machine learning and NLP.

III. TWITTER SENTIMENT ANALYSIS PROCEDURE

Before we perform the sentiment analysis on twitter data the data should be brought into proper form and sentiment relevant features need to be extracted. The steps followed in twitter sentiment analysis as shown in figure 1 are:

1. *Data collection*: Twitter allows researchers to collect tweets by using a Twitter API. One must have a twitter account to obtain twitter credentials (i.e. API key, API secret, Access token and Access token secret) which can be obtained from twitter developer site. Then install a twitter library to connect to the Twitter API.
 - a) “RT” is an acronym for retweet, which indicates that the user is repeating or reposting.
 - b) “#” stands for hashtag is used to filter tweets according to topics or categories.
 - c) “@user1” represents that a message is a reply to a user whose user name is “user1”.
 - d) Emoticons and colloquial expressions or slang languages are frequently used in tweets
 - e) External Web links (e.g. <http://amze.ly/8K4n0t>) are also frequently found in tweets to refer to some external sources.
 - f) Length: Tweets are limited to 140 characters.

2. Data Preprocessing

The data preprocessing can often have a significant impact on the performance of a supervised ML algorithm. The steps that are carried out in preprocessing of data are as follows-

- a) *Case Conversion*: All words are converted either into lower case or upper case in order to remove the difference between “Text” and “text” for further processing.

- b) Stop-words Removal: The commonly used words like a, an, the, has, have etc which carry no meaning i.e. do not help in determining the sentiment of text while analyzing should be removed from the input text.
- c) Punctuation Removal: Punctuation marks such as comma or colon often carry no meaning for the textual analysis hence they can be removed from input text.
- d) Stemming: Stemming usually refers to a simple process that chops off the ends of words to remove derivational affixes.
- e) Lemmatization: Deals with removal of inflectional endings only and to return the base or dictionary form of a word, which is known as the lemma.
- f) Spelling Correction: Spelling of the incorrect words can be corrected based on automated selection of more probable word.

3. Feature Extraction

Once the tweets are pre-processed we need to extract features relevant for sentiment analysis. Some of the features include:

- a) Term presence and frequency: It usually consists of n-grams of words and their frequency counts
- b) Parts of speech tagging: words in the text are tagged with their respective parts of speech in order to extract adjectives nouns verbs which add meaning to the sentiment.
- c) Opinion words and phrases: words or phrases that indicate opinion of the text
- d) Negation: presence of words like 'not', 'nor', 'neither' may reverse the sentiment of whole sentence. E.g. :“not good “
- e) Twitter specific features: Presence of emoticons in tweets, positive or negative hashtags are all twitter specific features which add meaning to the sentiment.

4. Training and Testing machine learning classifier:

After features are selected a Machine learning classifier is chosen for sentiment analysis. The classifiers are explained in section 3. Training data is used to train the classifier and its performance is measured using test data.

IV. MACHINE LEARNING ALGORITHMS FOR SENTIMENT CLASSIFICATION

A. Naïve Bayes classifier

The Naive Bayes classifier is the simplest (as the name suggests) and most commonly used classifier. Naive Bayes classifier works very well for text classification as it computes the posterior probability of a class, based on the distribution of the words (features) in the document. The

model uses the Bag of words feature extraction. It assumes that the features are independent of each other. It uses Bayes Theorem to predict the probability that a given feature[2]:

$$P(\text{label}/\text{features}) = P(\text{features}/\text{label}) * P(\text{label}) / P(\text{features})$$

$P(\text{label})$ is the prior probability of a label or the likelihood that a label is observed. Given a feature, $P(\text{features}|\text{label})$ is the prior probability that feature set is being classified as a label. $P(\text{features})$ is the prior probability that a given feature set is occurred. Given the Naive assumption which states that all features are independent of each other, the equation could be rewritten as follows:

$$P(\text{label}|\text{features}) = \frac{P(\text{label}) * P(f_1|\text{label}) * \dots * P(f_n|\text{label})}{P(\text{features})}$$

B. Support vector machines

The main principle of SVMs is to find out linear separators or hyperplane in the search space which can best separate the different classes. There can be several hyperplanes that separate the classes, but the one that is chosen is the hyperplane in which the normal distance of any of the data points is the largest, so that it depicts the maximum margin of separation.

Text classification are perfectly suited for SVMs because of the sparse nature of text, in which few features are irrelevant, but they tend to be correlated with one another and generally organized into linearly separable categories.[2]

C. Decision trees

Here, the training data space is represented in a hierarchical form in which a condition on the attribute value is used to partition the data. The condition on attribute values is the presence or absence of one or more words. The partition of the data space is done recursively until the leaf nodes contain certain minimum numbers of records which are used for the purpose of classification.

V. PERFORMANCE MEASURES

Once a classifier for sentiment analysis is selected, the trained model classifier must be validated using cross fold validation. The performance of the model can be determined using the following measures:

- a) *Accuracy*: It is measured by the fraction of number of correct predictions over total number of predictions. The accepted accuracy is usually in the range 70% to 90% . If a model is 100% accurate then it usually depicts that model overfits the data.
- b) *Precision*: This measure shows how accurately the model makes predictions w.r.t each class. It is measured by number of correct predictions over total number of true positives and true negative examples.
- c) *Recall*: This measure shows the completeness of the model w.r.t each class. It is measured by number of correct predictions over total number of true positives and false negative examples.
- d) *F-score*: It is measured as,

$$F\text{-score} = \frac{2 * \text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}}$$

VI. PROPOSED FRAMEWORK AND RESULT ANALYSIS.

The proposed approach follows the steps described in section 3 for sentiment analysis on twitter data. This work uses Apache Spark to obtain accurate results fast. Apache Spark is a quick and universally useful cluster computing system. It is popular for fast processing. Spark runs programs up to 100x speedier than Hadoop MR in memory, or 10x quicker on disk. Here, Naïve Bayes and decision trees machine learning algorithms are used for sentiment analysis. The two algorithms are compared based on the performance measures described in section 5.

This work performed analysis on datasets of different sizes and domains to demonstrate that the proposed framework works on data of all sizes and domains. Three different datasets of the following size and domains are considered for analysis:

Dataset Size: 200 tweets, 2000 tweets, 4000 tweets

Dataset Domains: IT Industry (Apple), Bank (ICICI), Telecom (BSNL)

Tweets ranging from the size 200 to 4000 were collected for all the domains mentioned above and the corresponding results were analyzed. 70% of each dataset is taken for training the classifier and remaining 30% is taken for testing. Different algorithms require different training and testing times. Decision tree and Naïve Bayes training and testing times are shown in table 1 below.

Algorithm	Dataset Size	Training Time (seconds)	Testing Time (seconds)
Multinomial Naïve Bayes	200	14.6713840961	3.004074096
	2000	10.5932500362	2.5033950805
	4000	7.55825281143	7.6055526733
Decision Tree	200	80.5690431595	0.161705970764
	2000	83.5256049633	0.199039936066
	4000	93.4236650467	1.77493691444

Table 1: Training and Testing time taken by Multinomial Naïve Bayes and Decision Tree for datasets of different sizes

Once the algorithms are trained, their performance is tested based on accuracy, precision, recall and F1-Score. The performance comparison of Multinomial Naïve Bayes and Decision Tree algorithms is shown in table 2.

Algorithms	Domain	Dataset Size (Tweets)	Accuracy In %	Precision In %	Recall In %	F1-Score In %
Multinomial Naïve Bayes	Apple	200	76	76.04	76.04	76.04
		2000	83.67	83.678	83.7	83
		4000	84.5775	84.6	84.68	84
	ICICI	200	78	78.4	78.4	78.4
		2000	82.7	82.678	82.7	82.5
		4000	84.8	84	84.6	84
	BSNL	200	77	77	77	77
		2000	83	83.2	83.24	83.2
		4000	85.6	85.6	85.6	85.6
Decision Tree	Apple	200	99	99	99	99
		2000	99.6	100	99	99.6
		4000	100	100	100	100
	ICICI	200	99.6	99.4	99.4	99.6
		2000	99.6	100	100	100
		4000	100	100	100	100
	BSNL	200	99	99	99	99
		2000	100	100	100	100
		4000	100	100	100	100

Table 2: Performance statistics of Algorithms based on dataset size and domains

The following points can be elaborated based on result analysis

- Multinomial Naïve Bayes does not perform as expected when supplied with small training dataset.
- Decision tree takes longer training time than Naïve Bayes.
- Decision tree takes very less time for predicting unseen data compared to multinomial Naïve Bayes.

- Decision tree performs better than Multinomial Naïve Bayes for datasets of varied sizes and domains.
- The proposed framework is domain independent and can operate on datasets of varied sizes.
- This work uses Apache Spark Cluster for processing. Hence, the proposed framework is able to produce the results of analysis quickly.

VII. CONCLUSION

Sentiment analysis can be performed using lexicon based approach, machine learning based approach or hybrid approach. Lexicon based approach faces a disadvantage that the strength of the sentiment classification depends on the size of the lexicon (dictionary). As the size of the lexicon increases this approach becomes more erroneous and time consuming. This paper explains in detail various steps for performing sentiment analysis on twitter data using machine learning algorithms. A machine learning classifier requires a labeled dataset which is divided into train and test set. Once an appropriate dataset is collected, the next step is to perform preprocessing on data (tweets) by using NLP based techniques, followed by feature extraction method in order to extract sentiment relevant features. Finally, a model is trained using machine learning classifiers like Naïve Bayes, Support Vector Machines or Decision trees and is tested on test data. The performance of the model can be measured in terms of accuracy, precision, recall and F-score.

The proposed framework performs sentiment analysis using Multinomial Naive Bayes and Decision tree algorithms. The results show that Decision tree performs extremely well showing 100% accuracy, precision, recall and F1-Score. The proposed text analytics framework is also real-time, fast, scalable, and reliable as we use Apache Spark framework.

References

- [1]. Mr. S. M. Vohra, 2 Prof. J. B. Teraiya, "A Comparative Study Of Sentiment Analysis Techniques", Journal Of Information, Knowledge

- And Research In Computer Engineering Issn: 0975 – 6760| Nov 12 To Oct 13 | Volume – 02, Issue – 02 Pg 313-317
- [2]. Walaa Medhat ,Ahmed Hassan, "Sentiment analysis algorithms and applications:A survey" Shams Engineering Journal (2014) 5, 1093–1113
- [3]. Alessia D’Andrea Fernando Ferri, "Approaches, Tools and Applications for Sentiment Analysis Implementation", International Journal of Computer Applications (0975 – 8887) Volume 125 – No.3, September 2015
- [4]. S. Milstein, A. Chowdhury, G. Hochmuth, B. Lorica, and R. Magoulas. "Twitter and the micro-messaging revolution: Communication, connections" An O’Reilly Radar Report. 54 pages, November 2008.
- [5]. Lei Zhang, Riddhiman Ghosh, Mohamed Dekhil, Meichun Hsu, Bing Liu, "Combining Lexicon-based and Learning-based Methods for Twitter Sentiment Analysis", Hewlett-Packard Development Company, L.P.2011
- [6]. Neethu M S, Rajashree R, "Sentiment Analysis in Twitter using Machine Learning Techniques", 2013 Fourth International Conference on Computing, Communications and Networking Technologies (ICCCNT).
- [7]. Min-Ling Zhang et.al, Feature selection for multi-label naive Bayes classification (2009),Elsevier, Information Science journal
- [8]. P.D Turney, "Thumbs up or thumbs down?semantic orientation applied to unsupervised classification of reviews "in proceedings of 4th annual meetings for computational linguistics, pp. 417-424,2002
- [9]. 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), pages 79{86, 2002}
- [10]. Eun Hee Ko, Diego Klabjan, "Semantic Properties of Customer Sentiment in Tweets," 2014 28th International Conference on Advanced Information Networking and Applications Workshops
- [11]. Bing Liu, Sentiment Analysis and Opinion Mining, Morgan and Claypool Publishers, May 2012.p.18-19,27-28,44-45,47,90-101.
- [12]. Sang-Hyun Cho and Hang-Bong Kang, "Text Sentiment Classification for SNS-based Marketing Using Domain Sentiment Dictionary", IEEE International Conference on Consumer Electronics (ICCE), p.717-718, 2012.
- [13]. Monisha Kanakaraj and Ram Mohana Reddy Guddeti, "NLP Based Sentiment Analysis on Twitter Data Using Ensemble Classifiers" 3rd International Conference on Signal Processing, Communication and Networking (ICSCN),2015
- [14]. Sunil B. Mane, Yashwant Sawant, Saif Kazi, Vaibhav Shinde, "Real Time Sentiment Analysis of Twitter Data Using Hadoop" International Journal of Computer Science and Information Technologies, Vol. 5 (3) , 2014, 3098 – 3100
- [15]. Gaurav D Rajurkar, Rajeshwari M Goudar, "A speedy data uploading approach for Twitter Trend And Sentiment Analysis using Hadoop", International Conference on Computing Communication Control and Automation,2015
- [16]. Malladihalli S Bhuvanl, Vinay D Rao, " Semantic Sentiment Analysis Using Context Specific Grammar" International Conference on Computing, Communication and Automation (ICCCA2015)