**COPY RIGHT**

# ELSEVIER
## SSRN

Title SENTIMENTAL ANALYSIS OF TWITTER DATA USING NLP AND MACHINE LEARNING TECHNIQUES

Paper Authors

**Kotha Jyoshna Priya,  Dr.G.Sai Chaitanya Kumar, D.Varun Prasad**

USE THIS BARCODE TO ACCESS YOUR ONLINE PAPER

To Secure Your Paper As Per UGC Guidelines We Are Providing A Electronic Bar Code

# SENTIMENTAL ANALYSIS OF TWITTER DATA USING NLP AND MACHINE LEARNING TECHNIQUES

[1] Kotha Jyoshna Priya, [2] Dr.G.Sai Chaitanya Kumar, [3] D.Varun Prasad

[1] PG Student, Department of CSE, DVR&Dr.HS MIC College of Technology, Kanchikacherla,AP

[2] Assoc. Professor, Department of CSE, DVR&Dr.HS MIC College of Technology, Kanchikacherla, AP,

[3] Assoc.Professor, Dept of CSE, DVR&Dr.HS MIC College of Technology, Kanchikacherla, AP

[1] jyoshna.kotta@gmail.com, [2] saichaitanyakumar@mictech.ac.in, [3] varunprasad@mictech.ac.in

**Abstract:**

Every social networking site, including Facebook, Twitter, Instagram, and others, has emerged as a major information source. It has been discovered that a corporate organisation can benefit from data extraction and analysis from social networking sites for their product promotion. One of the most common platforms for people to voice their opinions and sentiments about a product is Twitter. In our study, we examine how the general public feels about a product using data from Twitter. First, in order to filter tweets, we created a pre-processed data framework based on natural language processing (NLP). In order to analyse sentiment, we also use the Bag of Words (BoW) and Term Frequency-Inverse Document Frequency (TF-IDF) model concepts. This project aims to precisely categorise good and negative tweets using BoW and TF-IDF. The accuracy of sentiment analysis can be significantly increased by utilising TF-IDF vectorizer, and simulation results demonstrate the effectiveness of our suggested solution. Using NLP, we were able to analyse sentiment with an accuracy of 85%.

Index Terms—Natural language processing (NLP), Twitter, data mining, Sentiment analysis

**Introduction:**

The amount of data produced by internet services nowadays is enormous and is growing rapidly every day [1]. Microblogging is practised on social networking sites, where it has developed into a powerful instrument for communication among Internet users [2]–[3].Every business, large and small, has joined social networking sites to share their products and look for customer ratings. In order to measure customer happiness and improve their product, the corporation will employ sentiment analysis to understand how customers feel about their products. In particular, the established approach for sentiment analysis is used frequently to examine relationships between any device, famous person, sports team, and other entities.

After Facebook, Twitter is the second-largest social networking site, producing 21 million tweets each hour and 347,222 tweets per minute [1]. As a result, it opens the door for sentiment analysis and data

mining based on user tweets. Since sentiment analyses are a type of data mining, they allow for the observation of consumer sentiment toward a variety of topics and goods. It is also the foundation for techniques like machine learning, computational linguistics, biometrics, and natural language processing. Twitter is our platform of choice for sentiment analysis because it provides options for the sensitivity of articulated disposition. Twitter users can express their concise ideas via a short message because the character limit is 140 on the platform [4].

With the help of the Bag of Word (BoW) model and the TF-IDF (Term Frequency - In- verse Document Frequency) model concept, we have created an NLP-based pre-processed data framework to filter tweets in this paper. Tokenization, stemming, lemmatization, stop word removal, POS tagging, named entity recognition, coreference resolution, and text modelling like Bag of Word and TF IDF Model have all been accomplished using NLP techniques [5].When tweets are gathered using Twitter's streaming API, the major goal is to determine the emotion of each one by defining positive and negative polarity.

These tweets serve as basic data for us. At that moment, we apply the suggested technique that provides a tweet evaluation. Before making a purchase decision, the client will comprehend the feedback on the services from the sentiment analysis. The remainder of the essay is structured as follows. For analysing earlier work, Section II offers similar work. The suggested work's architectural overview and proposed algorithmic rule are briefly detailed in Section III. A basic description of sentiment analysis of a product is given in section IV. The simulated results and performance assessment of our suggested strategy are reported in section V. The paper is concluded in Section VI.

**Related Work:**

In [6], the authors suggested a machine learning approach for sentiment analysis using an existing twitter dataset. The idea of sentiment analysis utilising a suggested method that categorises Tweets automatically as good, negative, or neutral. Additionally, they are utilising a tree kernel and Part Of Speech (POS)-Specific polarity characteristics. In [7], an ontology-based approach for sentiment analysis is put out.

In order to identify the sentiment behind judgments and provide a description for such polarisation, the authors devised a model that combines domain ontology with natural language processing techniques. The technique tests were created employing the two distinct mediums of movies and digital cameras.[8].The techniques of sentiment analysis are trained in [9]–[10] to identify sentiment polarity, which allows them to automatically track down sentiments from various documents, blogs, sentences, or words. Through the integration of semantic technology, NLP, and information extraction, the authors of [11] created a new approach for semantic knowledge extraction from research materials and articles.

Using data cleansing, data extraction, and data consolidation, the authors of [12] proposed a novel method for extracting

structured data from emails. The supervised learning approach, according to the authors' theory in [13], is built on label datasets that are trained to provide useful outputs. Applying the Naive Bayes method, maximum entropy, and support vector machines to monitor the learning process enables successful sentiment analysis. In [14], the authors demonstrated that they could achieve an accuracy of up to 82.1 percent using the Naive Bayes algorithm. The accuracy of the sentiment analysis performed by the authors in [15] using the K-Nearest Neighbor classifier was 74.74 percent. The authors of [16] presented a survey on sentiment analysis of Twitter data using several methodologies. Different machine learning techniques, including Naive Bayes, Maximum Entropy, and Vector Machine Support, were employed to analyse sentiment and show the efficacy of various feature sizes.

**Proposed Work**

This section summarises our ideas and focuses on a method for conducting sentiment analysis on Twitter data. Figure 1 depicts the architectural overview defining a general process design for sentiment analysis. The created technique is based on three crucial components: Data Extraction from a specific project or product, Natural Language Toolkit (NLTK) preprocessing of the extracted Tweets, and Classifier model that determines the sentiment of each Tweet.
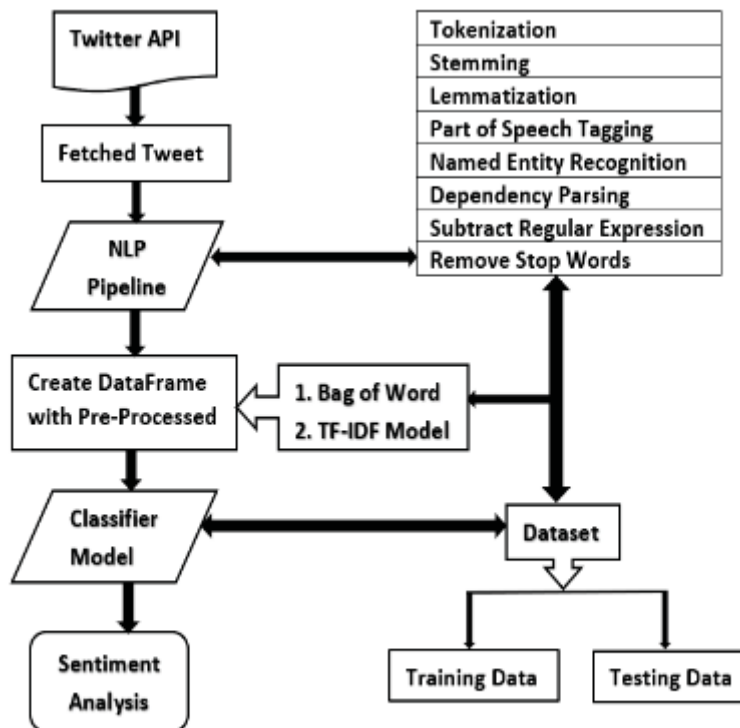


Fig. 1: Architectural overview of proposed work

Data extraction, tokenization, stemming, lemmatization, stopword removal, parts of speech tagging, named entity recognition, create a data frame, text modelling [17], and a classifier model are all used to perform sentiment analysis on the Twitter data. Each of these processes has its own algorithm and set of packages that must be

# International Journal for Innovative Engineering and Management Research
## A Peer Reviewed Open Access International Journal
www.ijiemr.org

processed. The TF-IDF model has been used to extract crucial terms from tweets in order to forecast sentiment. Pickle module, which has already completed the object serialisation process, aids in the construction of a classifier model.

In order to find the attitudes, we first constructed a classifier model using a dataset [18] that is displayed in Algorithm 1.Additionally, we created a sentiment analysis and tweet filtering pre-processed data framework based on natural language processing (NLP), as demonstrated in 2.In essence, we used the BoW model to classify documents and model their text before feeding them into our method for analysing tweets. BoW does not preserve any semantic information and accords all words the same significance. In order to identify the most crucial word in tweets, we override BoW and include the TF-IDF model. And it aids in highly accurate sentiment analysis. Our suggested system determines whether a tweet is favourable or negative based on the classifier model.

### Sentiment Analysis

This study presents the sentiment analysis approach we created for Twitter using Twitter data. We use the comparison between the two most popular devices, the iPhone and the Samsung, as an example. On an iPhone, we downloaded 100 tweets from Twitter, which are displayed. We have deleted stop words and superfluous information from tweets before applying sentiment analysis, which was covered in the previous section, using an NLP pipeline. We now employed "vectorizer" and "classifier," which are "clf" objects built from "tfidfmodel.pickle" and "classifier.pickle," to forecast sentiment. We display the sentiment forecast for each tweet for the iPhone device.

To test if our classifier model correctly predicts the sentiment of a tweet like the one in, we now take a sample of a tweet.

Since we do not analyse tweets based on emoticons, an emoticon was used in this tweet. The NLP pipeline removed the emoticon and many redundant words from the tweet, such as "https://t.co/c1yNjARvif," which was changed to "ynjarvif" by using the Python re package to preprocess the tweet. Since this word has no meaning, the TF-IDF model would return zero. Given 'zero' polarity, the classifier model in this instance predicts the sentence to be a negative statement.

It merely makes a deduction regarding the sentiment analysis. For better comprehension and visualisation, we have displayed several sentiments based on the amount of tweets collected, such as 50, 100, 500, and 1000, respectively. Within a year, two businesses released a large number of phones, some of which gained greater notoriety than those from competing companies. Therefore, it is impossible to state one company is more well-known than another because the difference in the proportion of favourable to negative tweets between the two phones isn't substantial. However, we can argue that both of these phone companies are more well-known than other phone companies.

### Natural Language Processing (NLP)

Natural Language Processing (NLP) is a cornerstone of modern AI and can be used to enable any number of different business use cases, from speech recognition to chatbots and entiment analysis. NLP can be used to help AI applications better understand what a given user wants and then as part of a larger platform, actually help the user to execute whatever action or operation they need. Jason Flaks understands what NLP is capable of more than most: He was the leading inventor of the conversational technology behind Microsoft's Kinect and Hololens.
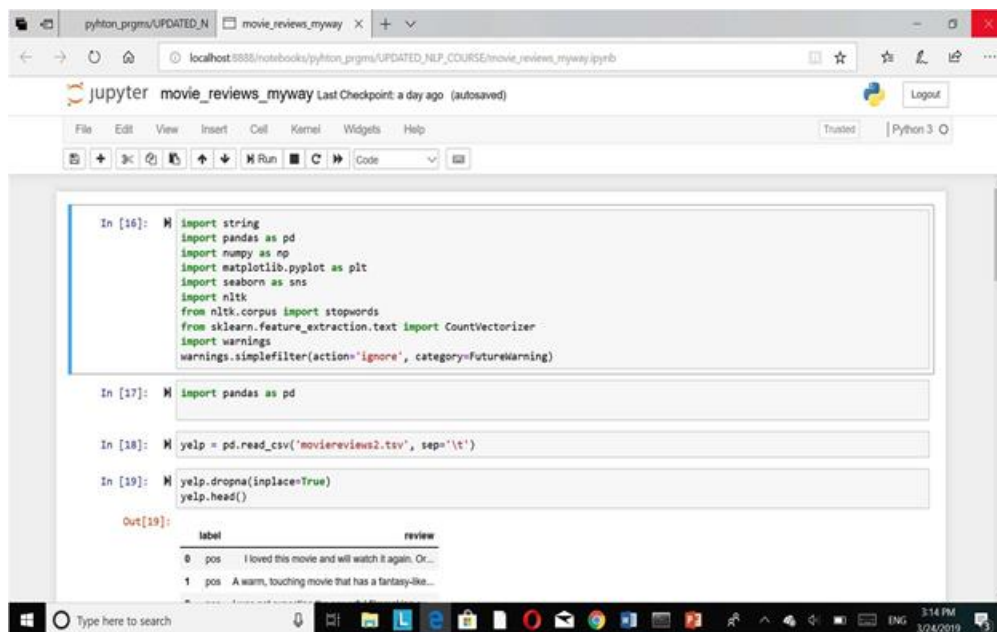
Now, Flaks is the CTO of startup Xembly, which is officially launching today. Flaks leads Xembly alongside CEO Pete Christothoulou, with the goal to use NLP and conversational AI to enable what the company refers to as an "AI chief of staff."

Helping business people and individuals to better organize and schedule their work lives is a complicated task. A top level business executive might be lucky enough to have a chief of staff to help schedule meetings, take notes and determine action items for follow up. Xembly aims to fill a similar role, with an automated NLP-powered AI platform. Modern conversational AI technologies and NLP can listen into a roomful of people and understand when more than just one person is talking and who is talking to whom. In Flaks' view, modern conversational AI technologies can actually be an active listener, even in a room full of people, rather than just waiting for a single voice to request a

single action. The real magic is, ultimately, what you can do with that," Flaks said. "I think that's where we continue to want to expand. The question is, how do we take that data and do meaningful things with it?"

The ability to transcribe notes is not a new thing for NLP. What Xembly is doing is going beyond just spoken dialogue to producing well-written prose for action items and meeting summaries. Flaks said that it takes a really complex set of machine learning models to be able to achieve that. He explained that Xembly runs its own machine learning models to detect action items from a conversation or interaction with Xena. The action items can include items like pulling out due dates and order requests. The Xembly system is also able to conduct topic segmentation for conversations to help users better understand the flow of a meeting.

**Results and Outputs**

## Conclusion:

In this paper, a basic yet novel methodology on sentimental analysis of motion picture surveys is performed utilizing promising supervised AI calculations. The out comes obtained concludes Logistic Regression as the best classifier among others in accomplishing 92% exactness for substantial number of movie reviews. Infuture, we attempt to explore its effectiveness considering enormous data set sutilizing the unsupervised and semi supervised AI methods. The NLP techniques for text pre processing are well performed to obtained data to for building the model. These techniques are far better than any other for the data preprocessing as it states that text classification is the main step forth emining of the data for the sentimental analysis. A conspicuous method to stretch out this work is add other order calculations to the, e.g., Conditional Random Fields or progressively expound gatherings. There are additionally a few highlights and highlight choice techniques that could be explored and a less innocent method for dealing with nullification. Instead of the straightforward treatment of in validation utilized here, away to deal with programmed enlistment of degree through a nullification finder could be utilized. The Future Sentimental Analysis work can be done on the group of the images to identify the emotions of the faces in the images and make the best filter upgrade for the group of images by identifying the good emotions the work can performed by the unsupervised algorithms

## Bibliography

[1] B. Pang and L. Lee, "Opinion mining and sentiment analysis," Foundations and trends in information retrieval,vol.2,pp.1- 135,2008.

[2] A.BifetandE.Frank,"Sentimentknowledgediscoveryintwitterstreamingdata,"in DiscoveryScience,2010,pp.1-15.

[3] B. Pang, L. Lee, and S. Vaithyanathan, "Thumbs up?: sentiment classification using machinelearning techniques," in Proceedings of the ACL-02 conference on Empirical methods in naturallanguageprocessing-Volume10,2002,pp.79-86.

[4] R.Basili,A.Moschitti,andM.T.Pazienza,"Languagesensitivetextclassification,"inRIAO,2000,pp.331-343.

[5] P. S. Jacobs, "Joining statistics with NLP for text categorization," in Proceedings of the thirdconferenceonAppliednaturallanguageprocessing,1992,pp.178-185.

[6] R. E. Schapire and Y. Singer, "BoosTexter: A boosting-based system for text categorization,"Machinelearning,vol.39,pp.135-168,2000.

[7] V.N.VapnikandV.Vapnik,Statisticallearningtheoryvol.1:WileyNewYork,1998.

[8] G.SaltonandM.J.McGill,"Introductiontomoderninformationretrieval,"1986.

[9] S. Dumais, J. Platt, D. Heckerman, andM. Sahami, "Inductive learning algorithms andrepresentations for text categorization," in Proceedings of the seventh international conference onInformationandknowledgemanagement, 1998,pp.148-155.

[10] S.M.Weiss,C.Apte,F.J.Damerau,D.E.Johnson,F.J.Oles,T.Goetz,etal.,"Maximizingtext-miningperformance,"IEEEIntelligentsystems,pp.63-69,1999.

[11] R. Feldman, "Techniques and

applications for sentiment analysis," Communications of theACM, vol. 56, pp. 82-89, 2013.

[12] S. Baccianella, A. Esuli, and F. Sebastiani, "SentiWordNet3.0:AnEnhancedLexicalRe sourceforSentimentAnalysisandOpinionMi ning,"inLREC,2010,pp.2200-2204.

[13] A. Pak and P. Paroubek, "Twitter as a Corpus for Sentiment Analysis and Opinion Mining,"inLREC,2010,pp.1320-1326.

[14] T.Blog, "Insightsintothe # World Cup conversation on Twitter," in TwitterBlog,ed, 2014.

[15] D. Terrana, A. Augello, and G. Pilato, "Automatic Unsupervised Polarity Detection on aTwitter Data Stream," in Semantic Computing (ICSC), 2014 IEEE International Conference on,2014,pp.128-134.

[16] L.Zhang,"SentimentanalysisonTwit terwithstockpriceandsignificantkeywordco rrelation,"2013.

[17] A.G.Jivani,"Acomparativestudyofs temmingalgorithms,"Int.J.Comp.Tech.App l,vol.2,pp.1930-1938,2011. . Hall, E. Frank, G. Holmes, B. Pfahringer, P. Reutemann, and I. H. Witten, "The WEKA datamining software: anupdate, " ACMSIGKDD explorationsnewsletter,vol.11,pp.10-18,2009.

[18] A. Genkin, D. D. Lewis, and D. Madigan, "Large-scale Bayesian logistic regression for text categorization,"Technometrics,vol.49,pp. 291-304,2007.

[19] H. Daumé III, "Notes on CG and LM-BFGS optimization of logistic regression," vol. 198, p.282,2004

[20] G.Sai Chaitanya Kumar, Dr.Reddi Kiran Kumar, Dr.G.Apparao Naidu, "Noise Removal in Microarray Images using Variational Mode Decomposition Technique" Telecommunication computing Electronics and Control ISSN 1693-6930 Volume 15, Number 4 (2017), pp. 1750-1756

[21] S. Gorintla, B. A. Kumar, B. S. Chanadana, N. R. Sai and G. S. C. Kumar, "Deep-Learning-Based Intelligent Pothole Eye+ Detection Pavement Distress Detection System," 2022 International Conference on Applied Artificial Intelligence and Computing (ICAAIC), 2022, pp. 1864-1869, doi: 10.1109 / ICAAIC53929.2022.9792696.