

Collective Intelligence Sentiment Analysis of Tweets using Machine Learning

Akshat Shrivastava, Anurag Sen, Amritansh Shrivastava, Sachin Singh, Nagesh Jadhav

Department of Computer Science and Engineering, Mit Adt University, India

Abstract— Microblogging websites like Twitter and Facebook, in this new era, is loaded with opinions and data. One of the most widely used micro-blogging site, Twitter, is where people share their ideas in the form of tweets and therefore it becomes one of the best sources for sentimental analysis. Opinions can be widely grouped into three categories good for positive, bad for negative and neutral and the process of analyzing differences of opinions and grouping them in all these categories is known as Sentiment Analysis. Data mining is basically used to uncover relevant information from web pages especially from the social networking sites. Merging data mining with other fields like text mining, NLP and computational intelligence we are able to classify tweets as good, bad or neutral. The main emphasis of this research is on the classification of emotions of tweets' data gathered from Twitter. This emotional opinion from the tweets can be used for analysis and decision making.

Index Terms— Machine Learning, Nlp, Sentimental Analysis, Social Media, Text blob, Twitter

1. INTRODUCTION

In the previous scarcely any years, there has been an immense development in the utilization of microblogging stages, for example, Twitter. By that development, organizations and media associations are progressively looking for approaches to utilize Twitter for data about what individuals think and feel about their items and administrations. Organizations, for example, Twitter (twitter.com), tweet feel, and Social Mention (www.socialmention.com) are only a rare sorts of people who publicize Twitter supposition investigation as one of their administrations.

While there has been a decent estimation of exploration on how slants are communicated in sorts, for example, online audits and news stories, how notions are communicated given the casual language and message-length imperatives of twitter has been considerably less contemplated.

Highlights, for example, programmed discourse labels and assets, for example, notion vocabularies have demonstrated helpful for supposition examination in different spaces, yet will they likewise demonstrate valuable for opinion investigation in Twitter? In this paper, we start to research this inquiry. [12]

Another test of microblogging is the amazing broadness of subject that is secured. It's anything but a misrepresentation to state that individuals tweet about everything without exception. In this manner, to have the option to fabricate frameworks to mine Twitter opinion about some random subject, we need a technique for

rapidly recognizing information that can be utilized for training. [12] [11] In this paper, we explore one method for building such data: using Twitter hashtags (e.g., best feeling, epic fail, news) to identify positive, negative, and neutral tweets to use for training three-way sentiment classifiers.

The online medium has become a significant way for people to express their opinions and with social media, there is an abundance of opinion information available. Using sentiment analysis, the polarity of opinions can be found, such as positive, negative, or neutral by analyzing the text of the opinion. Sentiment analysis has been useful for companies to get their customer's opinions on their products predicting outcomes of elections, and getting opinions from movie reviews. The information gained from sentiment analysis is useful for companies making future decisions. [10] Numerous customary methodologies in conclusion investigation utilizes the pack of words technique. The pack of words procedure doesn't think about language morphology, and it could inaccurately order two expressions of having a similar significance since it could have a similar sack of words. The connection between the assortment of words is considered rather than the connection between singular words. While deciding the general slant, the opinion of each word is resolved and consolidated utilizing a capacity. Pack of words likewise disregards word request, which prompts phrases with nullification in them to be inaccurately characterized

Twitter is an innovative service aired in 2006 with currently more than 550 million users [1]. The user created status messages are termed tweets by this service. The public timeline of twitter service displays tweets of all users worldwide and is an extensive source of real-time information. The original concept behind was to provide personal status updates. But the current scenario surprisingly witnesses tweets covering everything under the world, ranging from current political affairs to personal experiences. Movie reviews, travel experiences, current events etc. add to the list. Tweets (and in general) are different from reviews in their basic structure. [9] While reviews are characterized by formal text patterns and are summarized thoughts of authors, tweets are more casual and restricted to 140 characters of text. Tweets offer companies an additional avenue to gather feedback. Sentiment analysis to research products, movie reviews etc. aid customers in decision making before making a purchase or planning for a movie. Enterprises find this area useful to research public opinion of their company and products, or to analyze customer satisfaction. Organizations utilize this information to gather feedback about newly released products which supplements in improving further design. [8] Various methodologies which incorporate machine learning (ML) strategies, conclusion vocabularies, cross breed approaches and so on have been demonstrated valuable for feeling investigation on formal writings. Yet, their adequacy for extricating assessment in microblogging information should be investigated. A cautious examination of tweets uncovers that the 140-character length text confines the jargon which grants the estimation. The hyperlinks frequently present in these tweets thusly limit the jargon size. The differed spaces examined would most likely force obstacles for preparing. The recurrence of incorrect spellings and slang words in tweets (microblogs when all is said in done) is a lot higher than in other language assets which is another obstacle that should be survived. On the reverse way around the gigantic volume of information accessible from microblogging sites on shifted areas are exceptional with other information assets accessible. [8] [8] Microblogging language is portrayed by expressive accentuations which pass on a great deal of estimations. Intense lettered expressions, outcries, question marks, cited text and so forth leave scope for estimation extraction. The proposed work endeavors a novel methodology on twitter information by collecting an adjusted extremity vocabulary which has gained from item audits of the areas viable, the tweet explicit highlights

and unigrams to assemble a classifier model utilizing ML procedures.

2. LITREATURE REVIEW

Assessment examination is a developing territory of Natural Language Processing with research running from report level characterization (Pang and Lee 2008) to learning the extremity of words and expressions (e.g., (Hatzivassiloglou and McKeown 1997; Esuli and Sebastiani 2006)). Given the character confinements on tweets, ordering the notion of Twitter messages is generally like sentencelevel opinion investigation (e.g., (Yu and Hatzivassiloglou 2003; Kim and Hovy 2004)); be that as it may, the casual and specific language utilized in tweets, just as the very idea of the microblogging space make Twitter conclusion examination an altogether different errand. It's an open inquiry how well the highlights and methods utilized on increasingly all around framed information will move to the microblogging area. [7] Just in the previous year there have been various papers taking a gander at Twitter estimation and buzz (Jansen et al. 2009; Pak and Paroubek 2010; O'Connor et al. 2010; Tumasjan et al. 2010; Bifet and Frank 2010; Barbosa and Feng 2010; Davidov, Tsur, and Rappoport 2010). Different analysts have started to investigate the utilization of grammatical feature includes yet results stay blended. Highlights basic to microblogging (e.g., emoji's) are likewise normal, however there has been little examination concerning the handiness of existing slant assets created on no microblogging information. Analysts have additionally started to explore different methods of naturally gathering preparing information. A few analysts depend on emoji's for characterizing their preparation information (Pak and Paroubek 2010; Bifet and Frank 2010). (Barbosa and Feng 2010) misuse existing Twitter feeling destinations for gathering preparing information. (Davidov, Tsur, and Rappoport 2010) likewise use hashtags for making preparing information, yet they limit their trials to conclusion/non-notion grouping, instead of 3-way extremity order, as we do. Geetika and Divakar have examined the ML approaches for feeling investigation of Twitter information. The informational collection is preprocessed and the descriptors are extricated to get the element vector list. AI calculations like naive Bayes and SVM are applied to get the substance similitude [6] The drawback of this approach is that it uses the twitter data set which is labeled until that time. [6] The above strategy requires a great deal of manual work and a talented expert to remove and distinguish the opinion of the tweets. Further, the

investigation of twitter information must be done in close to ongoing with a deferral of 1-2 minutes for finding general feelings with high precision. [7] Twitter analysis is generally utilized for breaking down the client's recognitions as positive and negative by means of tweets. Informal communities have changed the method of dynamic by the clients and opinion examination calculations measure the view of clients computationally. The tweets present a gigantic volume of conclusion writings and data must be separated from these tweets in a convenient way. Dictionary based methodologies and ML based methodologies are utilized for assessment extraction from the tweets. These strategies require marked information to prepare the classifiers. Henceforth, this methodology is costly for new arrangement of information and the separated notion information is spoken to utilizing pie outline and html page which doesn't give an effective information representation. The content generated by the user provides a solid base for decision making in various fields like education, advertising, business intelligence and political polls. Analyzing the online available data and categorizing the public opinions from the highly unstructured twitter data is a challenging task. Supervised machine learning algorithms are proposed and their evaluation shows that their accuracy is moderate and performance is also low. Further, the techniques applied for analysis domain and language specific. Therefore, to overcome all the above stated issues, it is proposed to develop and implement a methodology for fetching real time tweets and perform sentiment analysis without skilled professional and complicated queries. The twitter data application also visualizes the live data stream and generates categorized output results. [5] These results decide the general feeling for any use of burst and stylish subjects effortlessly and high precision. Subsequently, it is conceivable to get supposition from the developing theme tweets in the genuine situation.

3. MATERIALS AND METHOD

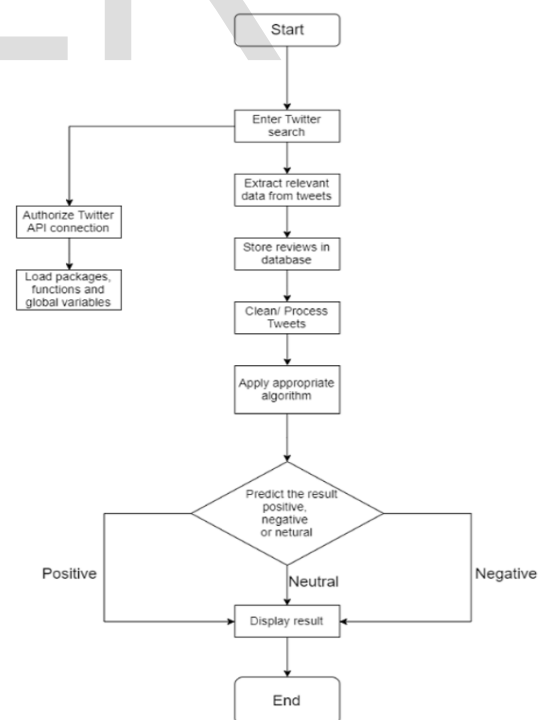
The Text Blob package for Python is a convenient way to do a lot of Natural Language Processing (NLP) tasks. For example: From text blob import Text Blob Text Blob ("not a very great calculation"). sentiment This tells us that the English phrase "not a very great calculation" has a polarity of about -0.3, meaning it is slightly negative, and a subjectivity of about 0.6, meaning it is fairly subjective. [4] There are helpful comments like this one, which gives us more information about the numbers we're interested.

In this technique we utilize text blob as a strategy to discover the extremity of the content (positive content, negative content or unbiased content). The tweets are imported from the Twitter utilizing the (API) gave by the Twitter Developer. From these API different fields like tweets, source, retweets, likes, language, client and so on can be rejected. In the wake of gathering these information, we can examinations the different celebrated individual considerations on an event or occasion. It clarifies the extraction of tweets id from twitter through API, at that point preprocess the information that are extricated. Preprocessing incorporates prohibition of undesirable fields, isolating the fields significant for investigation. When the fields are removed and isolated CSV is made. [4] [3] Using this CSV, the length of the message, Likes, retweets for the id is extracted and various results are derived. With the scraped tweets, classify the tweets whether positive or negative or neutral.

In this paper, we utilized python to actualize nostalgic examination. A few bundles have used including tweepy and text blob. We can introduce the necessary libraries by following orders:

- pip install tweepy
- pip install text blob

The second step is downloading the dictionary by running the following command: python -m textblob. download



corpora. bagheri2017sentiment The text blob is a python library for text processing and it uses NLTK for natural

Fig 1: Architecture Design

language processing. Corpora is a large and structured set of texts which we need for analyzing tweets.

4. OUR APPROACH

In our methodology we utilized the live twitter dataset and broke down it. We utilized the structure where the pre-processor is applied to the crude sentences which make it increasingly proper to comprehend. Further, the diverse ML strategies prepares the dataset with highlights and afterward the semantic investigation offers an enormous arrangement of equivalent words and comparability which gives the extremity of the substance. The total portrayal of the methodology has been depicted in next sub segments.

A. PRE-PROCESSING OF DATASETS

The tweets contain a ton of feelings about the information which are communicated in various manners by people. The twitters dataset utilized in this work is as of now marked. Named dataset has a negative and positive extremity and subsequently the investigation of the information turns out to be simple. The crude information having extremity is exceptionally powerless to irregularity and repetition. The nature of the information influences the outcomes and subsequently so as to improve the quality, the crude information is pre-handled. It manages the planning that expels the rehashed words and accentuations and improves the proficiency the information data.

B. FEATURE EXTRACTION

The improved dataset after pre-handling has a great deal of particular properties. The component extraction strategy, separates the viewpoint (descriptive word) from the dataset. Later this descriptor is utilized to show the positive and negative extremity in a sentence which is helpful for deciding the assessment of the people utilizing unigram model. Unigram model concentrates the descriptive word and isolates it. [4] It discards the preceding and successive word occurring with the adjective in the sentences. For above example, i.e. "painting Beautiful" through unigram model, only Beautiful is extracted from the sentence.

C. TRAINING AND CLASSIFICATION

Supervised learning is an important technique for solving classification problems. In this work too, we applied various supervised techniques to get the desired result for sentiment analysis. [4] In next few paragraphs we have briefly discussed about the three supervised techniques i.e. naive Bayes, maximum entropy and support vector machine followed by the semantic analysis which was used along with all three techniques to compute the

similarity. NAIVE BAYES It has been used because of its simplicity in both during training and classifying stage. It is a probabilistic classifier and can learn the pattern of examining a set of documents that has been categorized. It compares the contents with the list of words to classify the documents to their right category. $C^* = \text{argmax}_{c \in \{1, \dots, d\}} P(c|d)$ Class c^* is allocated to tweet d , where, f speaks to an element and $\text{in}(d)$ speaks to the include of highlight f found in tweet d . There is an aggregate of m highlights. Boundaries $P(c)$ and $P(f|c)$ are gotten through most extreme probability gauges which are increased by one for smoothing. Pre-handled information alongside separated element is given as contribution to preparing the classifier utilizing credulous Bayes. When the preparation is finished, during grouping it gives the extremity of the conclusions. For instance, for the audit remark "I am cheerful" it gives Positive extremity as result

MAXIMUM ENTROPY

entropy expands the entropy characterized on the contingent likelihood circulation. It even handles cover highlight and is same as calculated relapse which discovers conveyance over classes. It likewise follows certain component special case imperatives Where, c is the class, d is the tweet, and w is a weight vector. The weight vectors choose the centrality of a component in grouping. It follows the comparative procedures as Bayes, talked about above and gives the extremity of the opinions

SUPPORT VECTOR MACHINE

Support vector machine dissects the information, characterize the choice limits and uses the bits for calculation which are acted in input space. The information are two arrangements of vectors of size m each. At that point each datum spoke to as a vector is ordered in a specific class. [6] [5] Now the task is to find a margin between two classes that is far from any document. The distance defines the margin of the classifier, maximizing the margin reduces indecisive decisions. SVM also supports classification and regression which are useful for statistical learning theory and it helps recognizing the factors precisely, that needs to be taken into account, to understand it successfully.

SEMANTIC ANALYSIS

After the preparation and order we utilized semantic examination. Semantic examination is gotten from the WordNet database where each term is related with one another. This database is of English words which are connected together. On the off chance that two words are near one another, they are semantically comparable. All the more explicitly, we can decide equivalent word like closeness. We map terms and inspect their relationship in

the metaphysics. [4] The key errand is to utilize the put away archives that contain terms and afterward check the comparability with the words that the client utilizes in their sentences. In this manner it is useful to show the extremity of the notion for the clients. For instance, in the sentence "I am cheerful" the word "upbeat" being a modifier gets chose and is contrasted and the put away component vector for equivalent words. Let us expect 2 words; 'happy' and 'fulfilled' will in general be fundamentally the same as the word 'glad'. Presently after the semantic examination, 'happy' replaces 'cheerful' which gives a positive extremity.

5. RESULTS

We will initially introduce our outcomes for the goal/emotional and positive/negative characterizations. These outcomes go about as the initial step of our arrangement approach. We just utilize the short-recorded highlights for both of these outcomes. This implies for the goal/abstract order we have 5 highlights and for positive/negative grouping we have 3 highlights. For both of these outcomes we utilize the Naive Bayes order calculation, since that is the calculation we are utilizing in our real order approach at the initial step. we make a condition while detailing the aftereffects of extremity arrangement (which separates among positive and negative classes) that solitary emotional marked tweets are utilized to figure these outcomes. However, in case of final classification approach, any such condition is removed and basically both objectivity and polarity classifications are applied to all tweets regardless of whether they are labelled objective or subjective. In Twitter the various famous personalities tweet their thoughts on their opinion on an occasion. From their thoughts, importance of that occasion and the polarity of their tweet are analyzed Following shows the sample output of the program:

Information as crude tweets is gained by utilizing the Python library "tweepy" which gives a bundle to straightforward twitter gushing API. For our specific application we repeat through all the tweets in our example and spare the real content substance of the tweets in a document given that language of the twitter is client's record is determined to be English. The first content substance of the tweet is given under the word reference key "text" and the language of client's record is given under "lang".

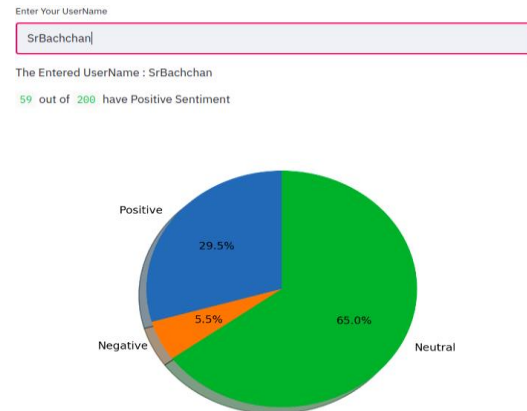


Fig 2: Analysis Pie Chart

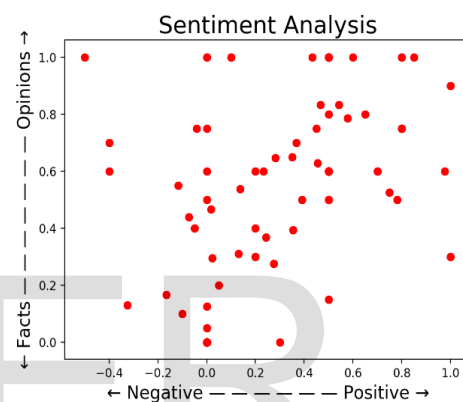


Fig 3: Analysis graph when keyword is used

1. 6. CONCLUSION

In this specialized paper, we talked about the significance of social network investigation and its applications in various zones. We concentrated on Twitter as and have executed the python program to execute wistful examination. We indicated the outcomes on various every day subjects. We understood that the unbiased sentiments are altogether high which appears there is a need to improve Twitter slant investigation.

Twitter analysis goes under the class of text and sentiment mining. It centres around breaking down the assumptions of the tweets and taking care of the information to an AI model to prepare it and afterward check its precision, with the goal that we can utilize this model for later use as per the outcomes. It includes steps like information assortment, text pre-handling, assumption identification, feeling grouping, preparing and testing the model. This research topic has evolved during the last decade with models reaching the efficiency of almost 85 [2] Be that as it may, it despite everything does not have the element of

decent variety in the information. Alongside this it has a great deal of utilization issues with the slang utilized and the short types of words. Numerous analysers don't perform well when the quantity of classes is expanded. [1] Also, it's still not tested that how accurate the model will be for topics other than the one in consideration. Hence sentiment analysis has a very bright scope of development in future.

7. REFERENCES

- [1] Apoorv Agarwal, Boyi Xie, Ilia Vovsha, Owen Rambow, and Rebecca Passonneau. Sentiment analysis of twitter data. In Proceedings of the Workshop on Language in Social Media (LSM2011), pages 30–38, 2011.
- [2] Oscar Araque, Ignacio Corcuera-Platas, J Fernando Sanchez-Rada, and Carlos an Iglesias. Enhancing deep learning sentiment analysis with ensemble techniques in social applications. *Expert Systems with Applications*, 77:236–246, 2017.
- [3] Hamid Bagheri and Md Johirul Islam. Twitter sentiment analysis. 2017.
- [4] Geetika Gautam and Divakar Yadav. Sentiment analysis of twitter data using machine learning approaches and semantic analysis. In 2014 Seventh International Conference on Contemporary Computing (IC3), pages 437–442. IEEE, 2014.
- [5] Ming Hao, Christian Rohrdantz, Halldór Janetzko, Umeshwar Dayal, Daniel A Keim, Lars-Erik Haug, and Mei-Chun Hsu. Visual sentiment analysis on twitter data streams. In 2011 IEEE Conference on Visual Analytics Science and Technology (VAST), pages 277–278. IEEE, 2011.
- [6] Zhao Jianqiang, Gui Xiaolin, and Zhang Xuejun. Deep convolution neural networks for twitter sentiment analysis. *IEEE Access*, 6:23253–23260, 2018.
- [7] Monisha Kanakaraj and Ram Mohana Reddy Guddeti. Performance analysis of ensemble methods on twitter sentiment analysis using nlp techniques. In Proceedings of the 2015 IEEE 9th International Conference on Semantic Computing (IEEE ICSC 2015), pages 169–170. IEEE, 2015.
- [8] Akshi Kumar and Arunima Jaiswal. Empirical study of twitter and tumblr for sentiment analysis using soft computing techniques. In Proceedings of the world congress on engineering and computer science, volume 1, pages 1–5, 2017.
- [9] MS Neethu and R Rajasree. Sentiment analysis in twitter using machine learning techniques. In 2013 Fourth International Conference on Computing, Communications and Networking Technologies (ICCCNT), pages 1–5. IEEE, 2013.
- [10] V Prakruthi, D Sindhu, and S Anupama Kumar. Real time sentiment analysis of twitter posts. In 2018 3rd International Conference on Computational Systems and Information Technology for Sustainable Solutions (CSITSS), pages 29–34. IEEE, 2018.
- [11] K Venkata Raju and M Sridhar. Sentimental analysis inclination a review. In 2017 International Conference on Current Trends in Computer, Electrical, Electronics and Communication (CTCEEC), pages 837–841. IEEE, 2017.
- [12] Varsha Sahayak, Vijaya Shete, and Apashabi Pathan. Sentiment analysis on twitter data. *International Journal of Innovative Research in Advanced Engineering (IJIRAE)*, 2(1):178–183, 2015.
- [13] M Trupthi, Suresh Pabboju, and G Narasimha. Sentiment analysis on twitter using streaming api. In 2017 IEEE 7th International Advance Computing Conference (IACC), pages 915–919. IEEE, 2017.