



# INTERNATIONAL JOURNAL OF RESEARCH IN SCIENCE & TECHNOLOGY

e-ISSN:2249-0604; p-ISSN: 2454-180X

## Leveraging Data Mining Tools and Techniques to Effectively Execute Sentiment Analysis on Social Media Data

Apoorva Khera

*Carmel Convent School, New Delhi*

**Paper Received:** 06<sup>th</sup> February, 2021; **Paper Accepted:** 20<sup>th</sup> March, 2021;

**Paper Published:** 29<sup>th</sup> March, 2021

DOI: <http://doi.org/10.37648/ijrst.v11i01.004>

### How to cite the article:

Apoorva Khera, Leveraging Data Mining  
Tools and Techniques to Effectively  
Execute Sentiment Analysis on Social  
Media Data, IJRST, January-March 2021,  
Vol 11, Issue 1, 34-40, DOI:  
<http://doi.org/10.37648/ijrst.v11i01.004>



## **ABSTRACT**

Lately, online social media has taken a fundamental part in communication and sharing of information. Countless users prefer social media as it is accessible to many individuals with no restrictions to share their understandings and instructive opportunities for growth and concerns through their status. Twitter API is handled to look for the tweets given the geo-area. Understudies' posts on the informal organization give us only worry about concluding the system-specific schooling systems growing experience. Assessing such information in an informal organization is a seriously difficult process. The proposed system will have a work process to mine the information, which incorporates both personal research and massive scope of information mining strategies. In light of the different noticeable subjects, we will classify tweets into various groups. A Naive Bayes classifier will be executed on searched data for personal investigation purposes to comprehend the information better. It involves a multi-name grouping strategy as each mark falls into various classifications, and every one of the characteristics is free. Will take name-based measures to break down the outcomes and distinguish them and the current sentiment analysis procedure.

## **INTRODUCTION**

Social sites like Twitter, WhatsApp and Facebook and so forth. Online Entertainment uses electronic and Internet devices to share views about data and encounters with humans more effectively. The Social Media Data Mining for Sentiment Analysis project is an online application for schools to their student's posted tweets.

The management for Students are Twitter presents to figure out issues in their informative experience. The authorities

given to the organization are Students posting their tweets via web-based social media to collect information connected with research growth opportunities. Students experience serious engagement on load, absence of social commitment and lack of sleep. The intricacy of students' encounters reflected in virtual entertainment content requires human Interpretation. To figure out what Student issues a tweet demonstrates is a more confusing task than deciding the opinion of a tweet, in any event, for a human-appointed authority. Subsequently, our review requires a subjective examination and is difficult to complete alone. The

proposed system needs to perform a subjective analysis using an order analysis rather than opinion mining. Sentiment analysis considers the client's perspective regarding a system or item and sorts it into neutral, pessimistic or good moods. In the proposed framework, looking through the data in light of the catchphrases, for example, engineer, understudies, grounds, class, teacher and lab in the Twitter data according to the geo-area, keyword and search id.

### **PROBLEM DEFINITION**

In the current system, the data includes synopses, meetings, polls, and homeroom exercises about the Student's instructive encounters and concerns. However, these conventional techniques are tedious and extremely limited in the hierarchy. The analysis doesn't check out in examining students' opportunities for growth which are enormous in volume with various Internet slang and the planning of the Student posting on the web. The sentiment analysis of the tweets doesn't cover a lot of important experiences since, in any event, for a judge to figure out what Student issues a tweet shows are a more confusing task than deciding simply the opinion of a tweet. These conventional techniques are extremely tedious and exceptionally limited in scale.

### **2. MODULE IN SYSTEM**

**Tweets Extraction and Preprocessing:** To extract tweets connected with the objective, we go through the entire dataset and release every one of the tweets containing the objective's keywords. Tweets are less formal and frequently composed of Adhoc than ordinary text records. Opinion mining techniques applied to raw tweets frequently perform ineffectively much of the time. This way, preprocessing techniques on tweets are fundamental for acquiring palatable outcomes on sentiment analysis: **Slang words interpretation:** Tweets frequently contain many shoptalk words (for example, Th, omg). These words are normally significant for sentiment analysis yet may not be remembered for opinion vocabularies. Since the sentiment analysis mechanism, we will utilize the opinion terminology. We convert these slang words into their standard structures utilizing the Internet Slang Word Dictionary<sup>1</sup> and afterwards add them to the tweets. **Non-English tweets are sifting:** Since the opinion examination techniques work for English texts, we eliminate all non-English tweets ahead of time. A tweet is viewed as non-English if more than 20% of its words (after interpretations of slang) don't show up in the GNU Aspell English Dictionary. **URL expulsion** A ton of clients remember URLs for their tweets. These URLs confound the sentiment analysis process.

We choose to eliminate them. Sensing Label Assignment To relegate opinion names for each tweet all the more without hesitation, we resort to two best-in-class opinion examination techniques. One is the SentiStrength3 mechanism [8]. This mechanism depends on the LIWC [10] sentiment dictionary. It works in an accompanying manner: first, dole out a feeling score to each word in the message as per the opinion vocabulary; then pick the most extreme good score and the greatest pessimistic score among those of generally individual words in the message; figure the amount of the most extreme good score and the most extreme pessimistic score, signified as Final Score; at last, utilize the indication of Final Score to demonstrate whether a tweet is good, neutral or pessimistic. This gathers the tweets from the data set to remove the tweets into the dataset. Organization Admin ControlPanel to visit then related understudy opportunity for growth to accessible in Twitter API question get.

### Naïve BAYES MULTI-LABEL CLASSIFIER

Naive Bayes is a straightforward probabilistic model in light of the Bayes rule with free element choice, which functioned admirably in the next order. This doesn't limit the number of classes or qualities to

manage. Asymptotically, Naive Bayes is the quickest learning analysis for the training stage. In this paper, we utilize the multinomial Naive Bayes model. In this formula,  $f$  addresses a component and  $n_i(d)$  address the count of element  $f$ ; found in tweet  $d$ ,  $m$  addresses the number of full highlights Considered. Class  $c^*$  is allowed to tweet  $d$ .

$$C^* = \arg \text{Max}_c P_{NB}(C|D)$$

$$P_{NB}(C/D) = \frac{(P(c) \sum_{i=1}^m P(f|c)^{n_i(d)})}{P(d)}$$

Boundaries  $P(c)$  and  $P(*)$  are obtained through the greatest probability gauges [11]. Bayes Rule Influence of one occasion's event on the likelihood of one more occasion is known as restrictive likelihood. From the likelihood hypothesis, the Bayes hypothesis considers contingent likelihood computation. Typically, the Bayes hypothesis will be utilized in information mining to choose substitute speculations. Bayes' hypothesis recipe for the contingent likelihood of A given B is as follows:  $P(A)$  is an earlier likelihood, and  $P(A|B)$  is the back likelihood of A given B.

$$\text{Polarity} = \frac{P(\text{positive Words}) / P(\text{Total Words})}{P(\text{Negative Words}) / P(\text{Total Words})}$$

Notwithstanding, this strategy turns out just for autonomous elements in light of the Standard English word reference and

neglects to catch question explicit opinions. Table-1 gives instances of positive, negative and unbiased tweets.

**Table 1: Tweet classification Example**

Sentiment	Tweet
Positive	I feel like I'm hidden from the world—life of an Engineering student.
Negative	I feel myself dying, #nervous.
Neutral	My problem is that I can never actually fall asleep without lying in bed for several hours.

**SYSTEM DESIGN**

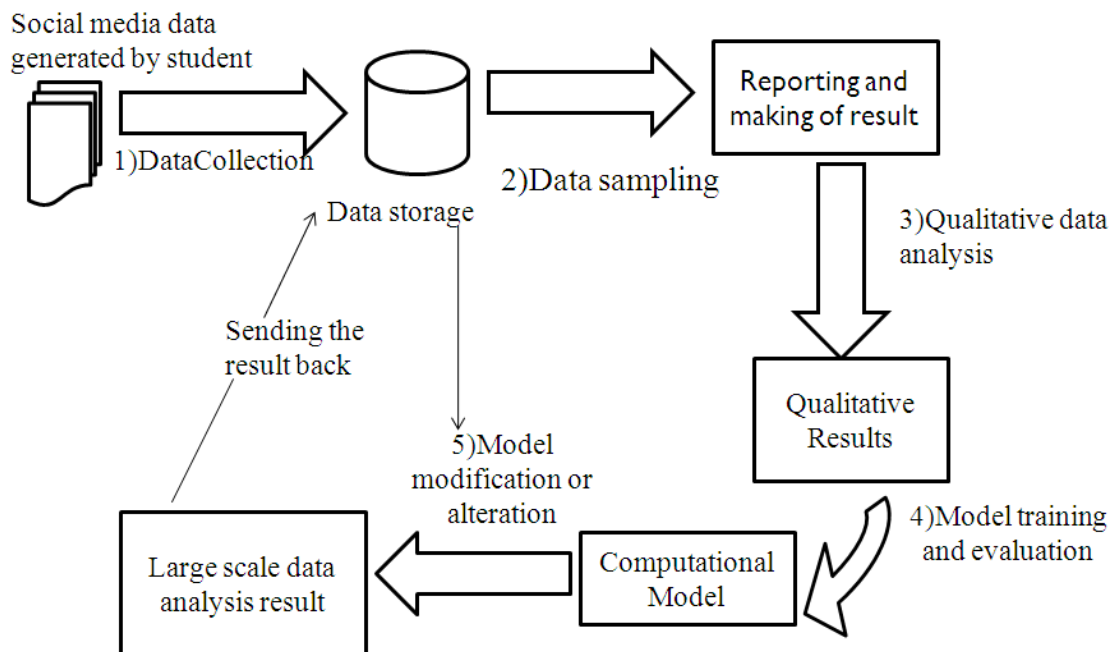


Fig. shows an example of a High-level system. To investigate public emotion analysis, There are two Latent Dirichlet Allocation (LDA) based models: (1) Foreground and Background LDA (FB-LDA) and (2) Reason Candidate and Background LDA (RCB-LDA). NaïveBayes, SVM, MaxEnt, ANN classifiers with highlights separated from Twitter information employing highlight extraction techniques, for example, unigram, Bigram and Hybrid (Unigram + Bigrams) for opinion examination. We perform information cleaning and standardization to eliminate stop words and focus keywords. We separate the objective-based dilated keywords model [7] by changing it and Twitter client information from the standardized data. We remove tweets connected with our fascinating targets (for example, "Understudy") and pre-process the removed tweets to make them more fitting for opinion investigation.

## CONCLUSION

We infer that informal community-based behaviour analysis boundaries can increase forecast accuracy. It is useful to professionals in learning analysis, instructive information mining and learning advancements. It gives a work process to breaking down virtual entertainment information for informational purposes that

defeats the significant restrictions of manual subjective analysis and the enormous scope of computational examination of client-produced literary substance. Our review can illuminate instructive directors, experts and other significant leaders to comprehend understudies' school encounters. Nonetheless, the presence of the relative multitude of elements fairly and equivalently is important to give precise outcomes. To understand the limitation that influences the outcomes, semantic keywords are also extremely valuable according to the perspective of the actual substance. Twitter-based informal communities give an incredible platform to evaluating popular assessments sensibly.

**Funding:** This research has no external funding.

**Conflicts of Interest:** Author declares no conflict of interest.

## REFERENCES

- [1] M. Rost, L. Barkhuus, H. Cramer, and B. Brown, "Representation and Communication: Challenges in Interpreting Large Social Media Datasets," Proc. Conf. Computer Supported Cooperative Workpp. 357-362, 2013.
- [2] J. Bollen, H. Mao, and X. Zeng, "Twitter mood predicts the stock market," J. Comput. Sci., vol. 2, no. 1, pp. 1-8, Mar. 2011.

[3] G. Siemens and P. Long, "Penetrating the fog: Analytics in learning and education," *Educause Review*, vol. 46, no. 5, pp. 30–32, 2011.

[4] R. Ferguson, "The state of learning analytics in 2012: A review and future challenges," Knowledge Media Institute, Technical Report KMI-2012-01, 2012

[5] G. Siemens and P. Long, "Penetrating the Fog: Analytics in Learning and Education," *Educause Rev.*, vol. 46, no. 5, pp. 30-32, 2011.

[6] M. Vorvoreanu, Q.M. Clark, and G.A. Boisvenue, "Online Identity Management Literacy for Engineering and Technology Students," *J. Online Eng. Education*, vol. 3, article 1, 2012.

[7] M.E. Hambrick, J.M. Simmons, G.P. Greenhalgh, and T.C.Greenwell, "Understanding Professional Athletes' Use of Twitter: A Content Analysis of Athlete Tweets," *Int'l J. Sport Comm.*, vol. 3, no. 4, pp. 454-471, 2010.

[8] B. O'Connor, R. Balasubramanian, B. R. Routledge, and N. A. Smith, "From tweets to polls: Linking text sentiment to public opinion time series," in *Proc. 4th Int. AAAI Conf. Weblogs SocialMedia*, Washington, DC, USA, 2010.

[9] 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.

[10] A. Abrahams, F. Hathout, A. Staubli and B. Padmanabhan, "Profit-Optimal Model and Target Size Selection with Variable Marginal Costs," 2013.

[11] R. Ferguson, "The state of learning analytics in 2012: A review and future challenges," Knowledge Media Institute, Technical Report KMI-2012-01, 2012.

[12] G. Siemens and P. Long, "Penetrating the fog: Analytics in learning and education," *Educause Review*, vol. 46, no. 5, pp. 30–32, 2011.

