

## Opinion Mining on Twitter Data of Movie Reviews using R

\*Ms. Md. Sania Sultana<sup>1</sup>, Mr. G. V Suresh<sup>2</sup>

<sup>1</sup>M. Tech in Dept. of Computer Science and Engineering, LBRCE College, Mylavaram, India.

<sup>2</sup>Associate professor in Dept. of Computer Science and Engineering, LBRCE College, Mylavaram, India.

Corresponding Author: Ms. Md. Sania Sultana

---

**Abstract:** Supposition investigation is essentially worried with the examination of feelings and assessments from the content. We can allude supposition examination as conclusion mining. Conclusion investigation finds and legitimizes the opinion of the individual regarding a given wellspring of substance. Web-based social networking contain a gigantic measure of the opinion information as tweets, websites, and updates on the status, posts, and so forth. Estimation examination of this to a great extent created information is exceptionally helpful to express the feeling of the mass. Twitter feeling investigation is dubious when contrasted with wide slant examination due to the slang words and incorrect spellings and rehashed characters. We realize that the greatest length of each tweet in Twitter is 140 characters. So it is essential to distinguish amend notion of each word. In our venture, we are proposing an exceedingly exact model of feeling investigation of tweets as for most recent  $r$  reviews of forthcoming Tollywood or Bollywood or Hollywood films utilizing Internet Movie Database (IMDb). With the assistance of highlight vector and classifiers, for example, Naïve Bayes, we are effectively characterizing these tweets as positive, negative and impartial to give feeling of each tweet.

**Keyword:** Naïve Bayes, Machine Learning, Twitter, Sentiment analysis, IMDb, Unigram.

---

Date of Submission: 07-07-2017

Date of acceptance: 28-07-2017

---

### I. Introduction

With the expansion in the notoriety of informal communication, small scale blogging and blogging sites, a colossal amount of information is produced. We realize that the web is the accumulation of systems. The age of the web has changed the way individuals express their musings and emotions. The general population are interfacing with each other with the assistance of the web through the blog entry, online discussion gatherings, and a great deal more. The general population check the audits or evaluations of the films before watching that motion picture in theaters. The amount of data is irrational for an ordinary individual to dissect with the assistance of credulous method .

Supposition examination is basically worried with the recognizable proof and grouping of sentiments or feelings of each tweet. Estimation examination is extensively arranged into the two sorts initial one is a component or angle based assessment investigation and the other is objectivity based feeling examination. The tweets identified with film surveys gone under the class of the element based supposition examination. Objectivity based conclusion examination does the investigation of the tweets which are identified with the feelings like despise, miss, love and so forth.

When all is said in done, different typical procedures and machine learning strategies are utilized to investigate the assessment from the twitter information. So in another way, we would say be able to that a supposition examination is a framework or model that takes the archives that broke down the information, and produces a point by point record compressing the sentiments of the given info report. In the initial step pre-preparing is finished. In the pre-handling we are evacuating the stop words, blank areas, rehashing words, emoticons, and #hash labels. To accurately arrange the tweets machine learning procedure utilizes the preparation information. Along these lines, this method does not require the database of words like utilized as a part of information based approach and hence, machine learning procedures are better and quicker.

The few techniques are utilized to separate the component from the source content. Highlight extraction is done in two stages: In the primary stage extraction of information identified with twitter is done i.e. twitters particular information is extricated. Presently by doing this, the tweet is changed into typical content. In the following stage, more elements are extricated and added to the element vector. Each tweet in the preparation information is related with a class name. This preparation information is passed to various classifiers and classifiers are prepared. At that point test tweets are given to the model and grouping is finished with the assistance of these prepared classifiers. So at long last we get the tweets which are characterized into the positive, negative and impartial.

## II. Literature Survey

There are two methods broadly used to recognize the assessments from the content. They are Symbolic methods and Machine Learning procedures [3].

### A. Sentiment analysis using Symbolic Techniques

A typical system utilizes the accessibility of lexical assets. Turney [4] recommended an approach for opinion investigation called 'sack of words'. In the specified approach, singular words are dismissed and just accumulations of words are considered. He accumulated word having descriptors or qualifier for the extremity of survey from a web crawler Altavista. A lexical database called WordNet [6] was utilized by Kamps et al [5] which decides a passionate matter in a word. WordNet conveys equivalent words and separation metric to discover the introduction of descriptors. To beat impediments in lexical substitution assignment, Baroni et al [7] built up a framework upheld by word space show formalism along these lines speaking to neighborhood words. EmotiNet adroitly spoke to the content that put away the structure of genuine occasions in a space. This was presented by Balahur et al [8].

### B. Sentiment analysis using Machine Learning Techniques

Under this system, there are two sets, to be specific a preparation set and a test set. By and large, the dataset which is gathered from various sources and whose conduct and yield esteems are known to us falls into the classification of preparing informational collections. Conversely with this, the datasets whose esteems or conduct are obscure to us are called as test informational indexes. Here various classifiers are prepared with preparing information and after that obscure information or we would say be able to a test information is given to this model to get coveted outcomes. Machine Learning comprises of different distinctive classifiers, for example, Ensemble classifier, k-eans, Artificial Neural Network and so on. These are utilized to characterize surveys [8]. Y.Mejova et al [1] in his exploration work suggested that we would use be able to the nearness of each character, recurrence of events of each character, word which is considered as nullification and so forth as components for making an element vector. He additionally demonstrates that we can successfully utilize unigram and bigram ways to deal with make include vector in Sentiment investigation. Domingos et al [10] proposed that Naive Bayes functions admirably for subordinate components for certain issue. Zhen Niu et al [11] found another model. This model depends on a Bayesian calculation. In this model, some effective methodologies are utilized for choosing the element, calculation of weight and grouping.

## III. Proposed Method

Different systems have been utilized to do supposition investigation of tweets. In our exploration, we have utilized the strategy for include vectors. The accompanying Figure demonstrates the whole proposed framework engineering. The proposed framework contains different periods of advancement. A dataset is made utilizing twitter posts of film audits. As we realize that tweets contain slang words and incorrect spelling. So we play out a sentence-level assessment investigation on tweets. This is done in three stages. In a first stage pre-preparing is finished. At that point Feature vector is made utilizing pertinent elements. At long last, utilizing diverse classifiers, tweets are arranged into positive, negative and unbiased classes.

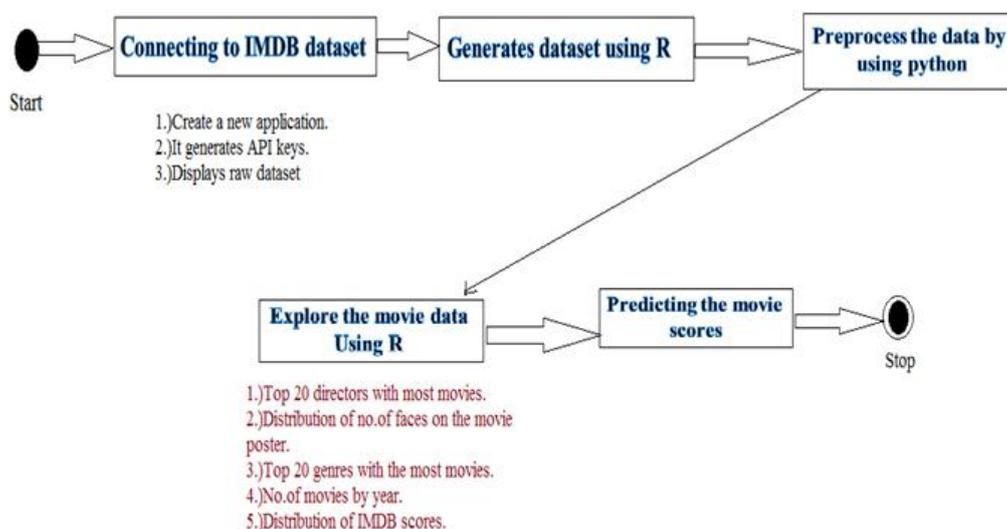


Fig 1: System model for our project with IMDb database.

### A. Creation of Dataset

A dataset is made utilizing twitter posts of motion picture surveys and related tweets about those films. The beneath table shows dataset utilized for preparing the classifiers and furthermore the tweets utilized for testing.

**Table I.** Statistics of the Dataset Used

Dataset	Positive	Negative	Neutral	Total
Training	600	600	600	1800
Testing	50	50	50	150

A dataset is created by taking 600 positives, 600 negatives, 600 neutral tweets.

### B. Pre-processing

The pre-preparing is the significant piece of our venture. In the pre-preparing, initial step is to change over the tweets into bring down case. Next are to maintain a strategic distance from the URL. Target name i.e. @username is supplanted by utilizing AT\_USER and hashtags are evacuated. Next, we have supplanted the rehashed character with the two events and expelled the blank areas.

### C. Sentiment Classification

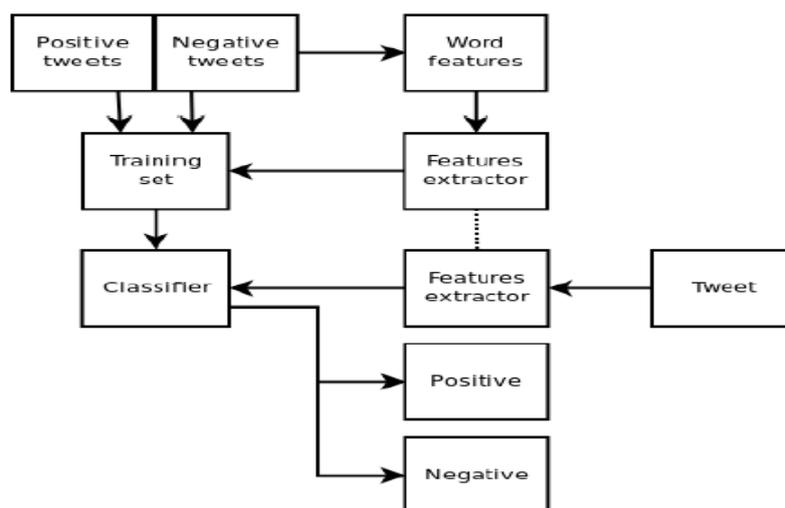
Subsequent to making a component vector, order is finished utilizing Naïve Bayes, Support Vector Machine and the execution is thought about

#### 1) Naïve Bayes Classifier

The primary favorable position of Naïve Bayes classifier is that it examinations each element autonomously. So it makes the utilization of the considerable number of components in the element vector. The Probability of Naïve Bayesian classifier is given as,

$$P\left(\frac{z}{b_j}\right) = \prod_{i=1}^m P(z_i/b_j)$$

Where the component vector is spoken to by  $z$  and  $b$  is the class mark (i.e. positive, negative, and nonpartisan). Another reason of utilizing Naïve Bayesian classifier is that it is easy to utilize and can be versatile. This classifier when contrasted with every single other classifier has high accuracy. Yet, the hindrance of it is that the precision and review are low. The gullible Bayesian classifier depends on the well known Bayes hypothesis in arithmetic.



**Fig.2** Sentiment Analysis Architecture

## IV. Details Of Experimentation And Discussion Of Results

This segment demonstrates the aftereffects of various tests that I executed amid my exploration. To begin with, the consequences of pre-preparing are appeared. These outcomes are trailed by consequences of the component determination and characterization. These consequences of the component determination and the order are joined since they can't be isolated from each other. The order calculation needs the components to group, and the element choice does not accomplish significant outcomes without the arrangement. The experimental details of the project are given as follows

## A. Datasets

The full preparing dataset contains the 21,000 tweets and put away in the CSV document. Out of these, we are utilizing 1200 tweets (600 positive tweets, 600 negative tweets, 600 unbiased tweets) for preparing the classifiers. These datasets are gathered from different sources and class names are physically commented on at whatever point class names are absent.

## B. Pre-processing of Tweets

In the initial step, the tweets are changed over to bring down case. So by doing this, we can get expressions of each tweet in a similar case (i.e. in bring down case). At that point in the following stage, every one of the URLs are disposed of and supplanted with typical content. At that point we have supplanted "@username" with bland word AT\_USER.

**Table II.** Example Showing Tweets and Feature Words

Positive Tweets	Feature Words
Bahubali2 The film is exceptionally positive. Celebrate Humanity. Doesn't take any religion or country's side.	'positive', 'Humanity', 'religion', 'country's', 'side'
Negative Tweets	Feature Words
AT_USER disappointed. Watched a movie. It is a waste of time.	'disappointed', 'watched', 'movie', 'waste', 'time'
I miss my mom and dad. I hate this life.	'miss', 'hate'
Natural Tweets	Feature Words
Not the best movie, but one time you can watch it	'Not', 'best', 'movie', 'but', 'time', 'one', 'watch'
By twitter, I am going to sleep now. Because tomorrow there is lots of work todo.	'sleep', 'now', 'tomorrow', 'lots', 'work', 'do'

## A. Classification of Tweets

### 1) Naïve Bayes Classifier

The arrangement utilizing Naïve Bayesian is done as takes after -

To start with, every one of the tweets and marks are passed to the classifier. In the subsequent stage, highlight extraction is finished. Presently, both these separated elements and tweets are passed to the Naïve Bayesian classifier. At that point prepare the classifier with this preparation information. At that point the classifier dump record opened in compose back mode and highlight words are put away in it alongside a classifier. After that the document is close.

## Pre-processing

**Table III:** Results of the pre-processing, applied cumulatively.

	English	non-English	unclassified
keyword filtering:	10,928	4,050	15,410
Twitter feature filtering:	14,529	7,242	8,617
overuse filtering:	14,629	7,568	8,191
replacement filtering:	14,710	7,546	8,132

In spite of the fact that the outcomes enhance, I expected that the substitution of acronyms and truncations would enhance the outcomes all the more altogether. This little change could be caused by the in total applying of the separating steps in light of the fact that a great deal of information is as of now sifted through. It could likewise be clarified by the method for arrangement of TextCat. Since the classifier is prepared on the general use of letters and letter mixes, acronyms and shortenings substitution don't change the appropriation essentially. An irregular example of 310 should cover 1% of the information since the dataset contains 30,388 occurrences. The physically named in-positions are contrasted and the anticipated names. The outcomes are appeared in the disarray grid in Table IV. With these numbers, different execution measures can be registered.

**Table IV:** Confusion matrix of a random sample of the results.

		Predicted class		
		English	Other	
Actual class	English	160	41	201
	Other	2	107	109
		162	148	

K-fold cross-validation is used to prevent overfitting and creating a more realistic view of the quality of the model.

**Table V:** Confusion matrix of model 1.1 according to the simple approach

		<i>Predicted class</i>			
		Positive	Neutral	Negative	
<i>Actual class</i>	Positive	224	274	72	570
	Neutral	459	1637	407	2,503
	Negative	126	293	235	654
		809	2,204	714	

However the circulation of grouped tweets is near the genuine dispersion of this dataset, it won't straightforwardly extrapolate to different datasets. This is on account of the arrangement precision on message level is somewhat low.

**Table VI:** Confusion matrix of model 1.2 with linear discriminant analysis based on strong positive, weak positive, weak negative and strong negative word frequencies.

		<i>Predicted class</i>			
		Positive	Neutral	Negative	
<i>Actual class</i>	Positive	40	519	11	570
	Neutral	46	2,392	65	2,503
	Negative	15	563	76	654
		101	3,474	152	

This is probably caused by the dominance of the neutral class in the dataset

**Table VII:** Confusion matrix model with the naive Bayes algorithm.

		<i>Predicted class</i>			
		Positive	Neutral	Negative	
<i>Actual class</i>	Positive	388	153	29	570
	Neutral	902	1,443	158	2,503
	Negative	222	194	238	654
		1,512	1,790	425	

For the third model, credulous Bayes show is utilized to analyze alternate models. The aftereffects of this arrangements are appeared in table VII. The arrangement precision is over 75.5%. Zooming in on the outcomes per class we see that credulous Bayes made a very unique forecast in contrast with the straightforward model. In this model, a considerable measure of nonpartisan cases are named positive. The high number of positive expectations could be caused by a more much of the time happening positive words than negative words in the component vector. When all is said in done, the innocent Bayes demonstrate predicts more cases to the positive and the negative classes. A clarification of this side effect is the premise of the element choice. The component choice is to be specific processed from the positive and negative tweets since impartial feeling words does not exist. In this manner the nonattendance of positive and negative conclusion words ought to demonstrate an unbiased tweet. Along these lines the quantity of nonpartisan expectations is diminished. The guileless Bayes calculation shows up rather to be decided for the negative and positive, caused by the freedom supposition.

## V. Conclusion

Consequently we reason that the machine learning system is extremely simpler and proficient than typical procedures. These systems are effectively connected to twitter notion investigation. Twitter conclusion investigation is troublesome in light of the fact that it is exceptionally difficult to distinguish enthusiastic words from tweets and furthermore because of the nearness of the rehashed characters, slang words, void areas, incorrect spellings and so on. Grouping exactness of the element vector is tried utilizing classifier like Nave Bayes. The presumption of Naïve Bayes that the information is free, turned out to be an amazing device in this examination. It was found by the creator that Machine learning calculations were more straightforward to actualize and more effective than different parts of the paper as they delivered a table which considered straightforwardness in the exactness of the Naive Bayes grouping. Generally speaking the half breed way to deal with opinion investigation considered an intensive examination of the information and performs well for a Twitter dataset. In any case, the precision of the Naïve Bayes classifier still leaves opportunity to get better this might be accomplished by better pre-preparing.

### Further development or research

The pertinence of slant investigation for future organizations and showcasing in utilizing watchwords and examination of the notions around that catchphrase by general society is just going to increment as the notoriety of Twitter becomes throughout the following couple of years. Be that as it may, as far as long haul improvement or research, the capacity of the twitter API to pull information that is more established, ought to be created and in addition another online networking API's so that estimation examination could be performed over some stretch of time, particularly in the domain of sociologies where specialists could enquire into social and political movements of sentiment on the web-based social networking locales. Similarly the absence of progress in feeling after some time on a few issues may be worth seeking after as a point of research for twitter slant examination. The convenience of such an opinion analyzer would take into account an intriguing examination of social and political issues.

### References

- [1] Kamps, J. Marx, M. Mokken, R. J. De Rijke, M. (2004 )"Using wordnet to measure semantic orientations of adjectives," Found in Neethu, M. Rajasree R.(2013) 'Sentiment Analysis in Twitter using Machine Learning Techniques' 4th International Conference on Computing, Communications and Networking Technologies (ICCCNT) Tiruchengode India. July 4-6 2013. IEEE.pp1-5 [Accessed April 25th 2014].
- [2] Y. Mejova, 'Sentiment analysis: An overview',mejova/publications/CompsYelenaMejova, vol. 2010-02-03, 2009, 2009.
- [3] E. Boy, P. Hens, K. Deschacht, and M. Moens, 'Automatic sentiment analysis in on-line text', 11th International Conference on Electronic Publishing, vol. 349360, 2007.
- [4] P. Turney, 'Thumbs Up or Thumbs Down? Semantic orientation applied to unsupervised classification of reviews', 40th annual meeting on association for computational linguistics, vol. 417424, 2002.
- [5] J. Kamps, M. Marx, R. Mokken and M. De Rijke, 'Using wordnet to measure semantic orientations of adjectives', 2004.
- [6] C. Fellbaum, 'Wordnet: An electronic lexical database (language, speech, and communication)', 1998.
- [7] D. Pucci, M. Baroni, F. Cutugno, and A. Lenci, 'Unsupervised lexical substitution with a word space model', Proceedings of EVALITA workshop, 11th Congress of Italian Association for Artificial Intelligence, Citeseer, 2009.
- [8] A. Balahur, J. Hermida and A. Montoyo, 'Building and Exploiting Emotional, a knowledge base for motion detection based on the appraisal theory model', Affective Computing, IEEE Transactions, vol. 3, 188101, 2012.
- [9] G. Vinodhini and R. Chandrasekaran, 'Sentiment analysis and opinion mining: A survey', International Journal, vol. 2, 6, 2012.
- [10] P. Domingos and M. Pazzani, 'On the optimality of the simple bayesian classifier under zero-one loss', Machine Learning, vol. 29, 2-3, 103130, 1997.
- [11] Z. Niu, Z. Yin and X. Kong, 'Sentiment classification for microblog by machine learning', Computational and Information Sciences (ICCIS), 2012 Fourth International Conference on, pp. 286-289, IEEE, vol. 286289, 2012.
- [12] L. Barbosa and J. Feng, 'Robust Sentiment Detection on Twitter from Biased and Noisy data', 23rd International Conference on Computational Linguistics: Posters, vol. 3644,, 2010.
- [13] A. Celikyilmaz, D. Hakkani-Tur and J. Feng, 'Probabilistic Model-Based Sentiment Analysis of Twitter Messages', Spoken Language Technology Workshop (SLT), 2010 IEEE, vol. 7984, 2010.
- [14] Y. Wu and F. Ren, 'Learning sentimental influence in twitter', Future Computer Sciences and Application (ICFCSA), 2011 International Conference, IEEE, vol. 119122,, 2011.
- [15] A. Pak and P. Paroubek, 'Twitter as a Corpus for Sentiment Analysis and Opinion mining', Proceedings of LREC, vol. 2010, 2010.
- [16] V. Peddinti and P. Chintalapoodi, V.M.Kiran, 'Domain adaptation in sentiment analysis of twitter', Analyzing Microtext Workshop, AAAI, 2011.

IOSR Journal of Computer Engineering (IOSR-JCE) is UGC approved Journal with Sl. No. 5019, Journal no. 49102.

Ms. Md. Sania Sultana. "Opinion Mining on Twitter Data of Movie Reviews using R." IOSR Journal of Computer Engineering (IOSR-JCE) 19.4 (2017): 19-24.