

# Sentiment Analysis of Movie Review using data Analytics Techniques

H. SWATHI<sup>1</sup>, S. S. ARAVINTH<sup>2</sup>, V. NIVETHITHA<sup>3</sup>, T. SARANYA<sup>4</sup>, R. NIVETHANANDHINI<sup>5</sup>  
<sup>1,2,3,4,5</sup> Dept. of CSE, Dhirajlal Gandhi College Of Technology, Salem, Tamilnadu

*Abstract -- Social media, and then the online platforms contain a huge amount of the data in the form of tweets, blogs, and updates on the status, posts, etc. Sentiment analysis of the data is handy to express the opinion of the group or any individual. Sentiment analysis is the analysis of emotions and opinions from any type of text. And the comparison of the Times of India movie reviews and Tweets of the same movies are taken using Machine Learning Algorithms. The primary motivation was to explore disparate data sources that had not encountered before that are also often used in data science, specifically Twitter and sentiment analysis. The main goal was to gain insight on market research analysis on social media, and traditional media sources. The resulting analysis is helpful to moviegoers and the movie industry in understanding the perception of reviews.*

**Indexed Terms:** Data Analytics, Movie Review, Twitter, SVM

## I. INTRODUCTION

Our primary motivation is to explore different data sources that we had not encountered before that are also often used in data science, specifically Twitter and sentiment analysis. We visualized a sentiment analysis comparison of New York Times movie reviews and Tweets of the same movies. We able to answer how favorably the New York Times reviewed a given movie compared to the average sentiment on Twitter. Our main goal was to gain insight on research analysis on social media, and traditional media sources. We wanted our resulting analysis to be useful to moviegoers, and the movie industry in understanding the perception of reviews.

## II. PROJECT DESCRIPTION

Sentiment analysis is a sub-domain of opinion mining where the analysis is focused on the extraction of emotions and opinions of the people towards a particular topic from a structured, semi-structured or unstructured textual data. In this project, we try to focus our task of sentiment analysis on IMDB movie

review database. We examine the sentiment expression to classify the polarity of the movie review on a scale of 0 as highly disliked to 4 as highly liked and perform feature extraction and ranking and use these features to train our multi-label classifier to classify the movie review into its correct label. There is a lack of strong grammatical structures in movie reviews, which follow the informal jargon, so we used an approach based on structured N-grams. On top of that, the comparative study on different classification approaches has been performed to determine the most suitable classifier to suit the problem domain. We finalized that our projected approach to sentiment classification supplements the current rating film rating systems used across the online and it will function as base to future researches in this domain. "Our approach using classification techniques has the good accuracy of 88.96%".

The Several multi class classification algorithms have been evaluated with given training data to find out the best algorithm for the task. While evaluating the accuracy of algorithms, same training set cannot be used as model may over fit to training data but it cannot predict anything useful for unseen data. To avoid this problem, we can use common practice to divide given training data to train set and test set. There are many approaches to divide given training data set to train set and test set. Initially hold out approach was used where 60% of original data set is used for training and remaining amount is used for testing. But there is still a risk of over fitting on the test set because the parameters can be strained until the algorithm performs optimally. The knowledge about the test set may leak into the model and evaluation metrics no longer report on generalization performance. Another part of the data set can be held out as validation set to solve this issue. So the work flow for evaluating is training proceeds on the training set, evaluation is done on the validation set,

and when the experiment seems to be successful, final evaluation can be done on the test set.

### III. LITERATURE SURVEY

#### 1. Sentiment Analysis of Movie Reviews: A Study on Feature Selection & Classification Algorithms:

The present era of net has become enormous Cyber information that hosts large quantity of information that is formed and consumed by the users. The information has been growing at associate exponential rate giving rise to a replacement business crammed with it, during which users specific their opinions across channels like Facebook, Twitter, Rotten Tomatoes, and Foursquare. Opinions that area unit being expressed within the type of reviews offer a chance for brand spanking new explorations to seek out collective likes and dislikes of cyber community. One such domain of reviews is that the domain of picture reviews that affect everybody from audience, film critics to the assembly company. The picture show reviews being denoted on the websites don't seem to be formal reviews however square measure rather terribly informal and square measure unstructured kind of descriptive linguistics. An opinion expressed in picture show reviews provides a terribly true reflection of the feeling that is being sent. The presence of such a good North American nation of sentiment words to precise the review impressed us to plot Associate in Nursing approach to classify the polarity of the picture show exploitation these sentiment words.

Sentiment Analysis could be a technology, which will be important within the next few years. With opinion mining, we are able to distinguish poor content from prime quality content. With the technologies on the market, we will able to apprehend if a pic has additional smart opinions than dangerous opinions and realize the explanations why those opinions are positive or negative. Much of the early research in this field was centered around product reviews, such as reviews on different products on Amazon.com, defining sentiments as positive, negative, or neutral. Most sentiment analysis studies are now focused on social media sources such as IMDB, Twitter and Facebook, requiring the approaches being tailored to

serve the rising demand of opinions in the form of text. Furthermore, playacting the phrase-level analysis of pic reviews proves to be a difficult task.

In this paper, we follow a lexical approach using the SentiWordNet to determine the overall polarity of the movie review. We analyze and study the features that affect the sentiment score of the movie review text. Also, we use the state of the art classification algorithms for the evaluation of performance and accuracy of the approach used. Also, we not only study the approach but try to have a deeper understanding of the problem domain.

#### 2. Sentiment Analysis of Movie Reviews Using Machine Learning Techniques:

Sentiment analysis is essentially involved with analysis of emotions and opinions from text. We can refer sentiment analysis as opinion drilling. Sentiment analysis finds and justifies the sentiment of the person with relation to a given supply of content. Social media contain vast quantity of the sentiment knowledge within the type of tweets, blogs, and updates on the standing, posts, etc. Sentiment analysis of this mostly generated knowledge is incredibly helpful to specific the opinion of the mass. Twitter sentiment analysis is hard as compared to broad sentiment analysis attributable to the slang words and misspellings and recurrent characters. We know that the utmost length of every tweet on Twitter is a hundred, and forty characters. So it's vital to spot correct sentiment of every word. In our project we tend to area unit proposing an extremely correct model of sentiment analysis of tweets with relation to the latest reviews of coming screen land or Hollywood movies. With the assistance of feature vector and classifiers like Support vector machine and Naïve mathematician, we tend to area unit properly classifying these tweets as positive, negative and neutral to administer sentiment of each tweet.

### IV. EXISTING SYSTEM

In this problem, it's common follow to divide given coaching knowledge to coach set and take a look at set. There square measure numerous approaches to divide given coaching knowledge set to coach set and take a look at set. Initially hold out approach was

used wherever hr of original coaching knowledge set is employed for coaching and remaining quantity is employed for testing. But there's still a risk of over fitting on the take a look at set as a result of the parameters are often pinched till the algorithmic rule performs optimally. This way, data concerning the take a look at set could leak into the model and analysis metrics now not report on generalization performance. Another a part of the data set is often command out as validation set to unravel that issue.

So the work flow for evaluating is training proceeds on the training set, evaluation is done on the validation set, and when the experiment seems to be successful, final evaluation can be done on the test set.

➤ Disadvantages:

- The system needs to be hosted on cloud to receive and process country wide result.
- The writing can misguide you about the reality.
- Sometimes they serve as spoilers ruining the experience for you.Nothing would seem new and fresh.
- The system must be given proper inputs otherwise system can produce wrong result.

## V. PROPOSED SYSTEM

1. There is a lot of scope for doing more work on review in future. Some of them will be listed below
2. Various techniques have been used to do sentimental analysis of tweets. In this research method of feature vectors is used
3. Getting feedback from users to improve the user interface.
4. Using text analytics to find out the opinions of common people.
5. Using more new advanced algorithms for increased accuracy in predictions.
6. The tweets contain slang words and misspelling, so perform a sentence level sentimental analysis on tweets
7. K-Nearest Neighbor, Naïve Bayes and Random Forest featured algorithms used

➤ Advantages:

- The major advantage of support vector machines is effectiveness in high dimensional areas.
- Also it uses a subset of training points in the decision function called support vectors, so it is also memory efficient.
- SVM is when training data is highly unbalanced, resulting model tends to perform well on majority data but perform bad on minority data. In sickest library, different kinds of kernels like linear, rbf and polynomial square measure are provided.

## VI. SYSTEM ARCHITECTURE

1. We began by identifying our project goals and scope, assessing potential data sources, and retrieving the data.
2. We first performed an API pull of a list of recent New York Times movie reviews.
3. We then performed an API pull of recent tweets. We learned that Twitter provides only a limited amount of search history for publicly available tweets. This resulted in associate adjustment of our project scope.
4. We then adjusted the date range of our New York Times movie list to contain only the movies with available Twitter data.
5. We ran a sentiment analysis on the text of the tweets and the text of the movie reviews for the given set of movies.;
6. We created data frames for our data sources, erased and transformed the data, and performed analysis on the resulting datasets.

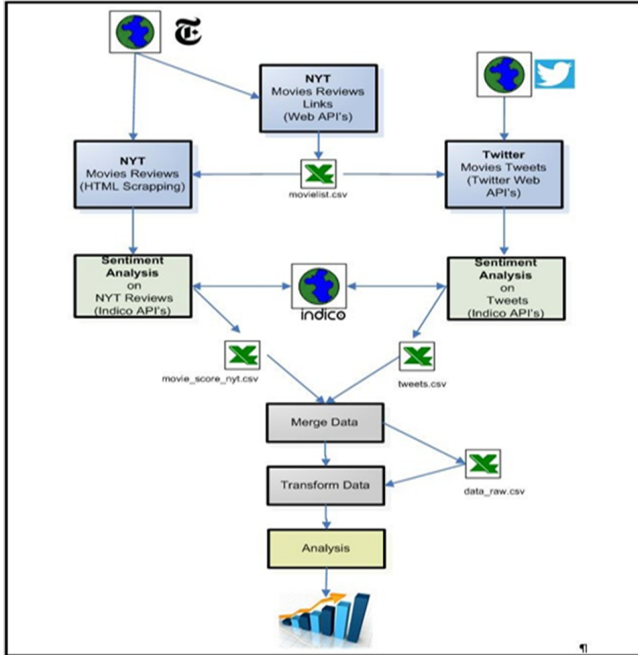


Fig. 1: System Architecture for sentiment Analysis

VII. RESULTS

1. Twitter Sentiment Analysis:

- Once we tend to establish the list of available New York Times movie reviews, we ran the Twitter search using the movie title as a search term. To get best Twitter search results for movies, we added the word “movie” to our search term.
- This chunk retrieved the Twitter search results for a term. searchTwitter has some other settings that could probably improve results.
- We used the indico package to perform sentiment analysis, resulting in a score of 0-1 for each tweet, 1 being extremely positive.
- We decided to use unique tweets only, therefore eliminating any retweet data.
- For the benefit of our analysis, we filtered the Twitter search results to pull only the tweets which were coded with geographical location.
- We chose the following cities at random and pulled tweets for only these cities.

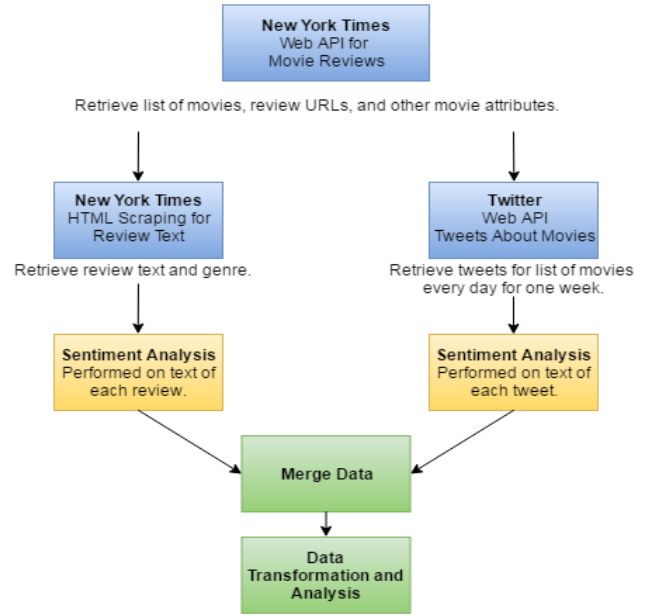


Fig. 2: Workflow of Movie Review Analysis

2. New York Times Sentiment Analysis Litmus Test:

We performed a sentiment analysis litmus test for New York Times movie reviews. The Times does not offer “stars” to rate their movie reviews however favorable reviews are labeled with a “Critics’ Pick.” As a check for the sentiment analysis accuracy, we analyzed the results as follows:

We considered a “good” score  $\geq 0.75$  and expected the movie to be tagged with “Critics’ Pick” (value = 1).

We looked at Pick/Good (True Positive), Pick/Bad (False Negative), Not Pick/Good (False Positive), Not Pick/Bad (True Negative).

- Percentage of False Positive: 35%
- Percentage of False Negative: 5%
- Percentage of True Positive: 15%
- Percentage of True Negative: 45%

The number of false negative were indicators of however troublesome it had been to accurately live the key and sentiment in an exceedingly document. There were 2 movies that fell into the class of false

negative. The movies were counseled by the New York Times critic; but, the sentiment analysis score was less than the threshold we considered. By reading the reviews, for these movies, “Viva” and “A Hologram for a King,” the sentiment analysis results were understandable. Both movies were dramas with components of the review portraying troublesome, dark, and negative terms. The false positive results were more easily explained. The the New York Times gave neutral to positive reviews while not recommending the films. Also, a score of zero.75 could be too low a limit to be expected to be a recommendation. We thus thought-about scores below .50 as negative, scores between .50 and .80 as neutral, and above .80 as positive.

We reran the analysis based on these new limits. False negatives were scores  $< 0.50$  but with a recommendation; false positives were scores  $\geq 0.80$  without a recommendation, and false neutrals were scores within interval (.50-.80) with a recommendation.

#### VIII. CONCLUSION

1. The paper mainly addresses the comparative study of different machine learning algorithms that can be used to extract sentiments from text
2. Naïve bayes and SVM with the best accuracy can be considered as a benchmark for all other algorithm
3. Intelligent system can be developed which can provide the user with comprehensive reviews of movies, products, services etc.
4. This provides information that could help to improve the predictions in further research work
5. This is simple and efficient.
6. It can be concluded that cleaner data, better the performance of an algorithm in predicting the success rate of the movies.

#### IX. ACKNOWLEDGMENT

Our sincere thanks to the staffs that helps us for the preparation of this paper about the project.

#### REFERENCES

- [1] Tirath Prasad Sahu, Sanjeev Ahuja, 'Sentiment Analysis Of Movie Reviews: A Study on Feature Selection And Classification Algorithm', 'IEEE-2016'.
- [2] Palak Baid, Apoorva Gupta, et. al, 'Sentiment Analysis Of Movie Reviews using Machine Learning Techniques', 'International journal of Computer Applications'.
- [3] Kuat Yessenov, Sasa Misailovic, 'Sentiment Analysis of Movie Review Comments', 2009.
- [4] Dorothy Aku Allotey, Regina, Saskatchewan, 'Sentiment Analysis and classification of online reviews using Categorical Proportional Difference' 2011.
- [5] Ieva Stalliunaite, Ben Bonfil, 'Breaking Sentiment Analysis of Movie Review', 2017.