

# Sentiment Identification using Movie Reviews

Devika Sharma, Rajesh Babu

*Abstract: Surveys are a noteworthy basis for the appraisal of the item, however the normal human pursuer can just investigate a limited number of the, under the given time imperatives which is the reason we need a computerized framework to break down these audits and think of a yield depicting the slant behind these surveys. The purpose of this project is to utilize a lot of heterogeneous highlights (Machine learning based and dictionary or lexicon based highlights) and administered learning calculations, for example, Naïve Baye's classifier and LSVMs to precisely examine an item dependent on the opinions of the surveys created for it.*

## I. INTRODUCTION

People's decisions are regularly vigorously one-sided by the individuals who encompass them which is the reason audits are a key element for examination of an item whether its motion pictures, magnificence items, garments. In certainty even the significant choices that we take in life are frequently checked on by higher specialists before we settle them. The paper utilizes ideas of Data mining and Natural Language Processing with the goal that human pursuers can concentrate and extract the given reviews in lesser time. It inputs a huge number of surveys and utilizing the idea of notion examination it audits the item for the client.

The purpose of this project can be broadly divided into two parts:

- To concoct an effective technique for Sentiment Analysis
- To analyse the given reviews for identifying the sentiment associated with them.

The paper follows a modular programming approach which is further explained in the method section. All the modules are executed one after the other in a pipeline architecture.

[1] "Sentiment Analysis of Movie Reviews using Heterogeneous features":

This paper aims at increasing the efficiency of the current model by increasing the accuracy of the results generated. This is done by incorporating a bi-feature method which uses two different types of features:

Lexicon based and Machine Learning based features. This helps in increasing the number of true positives and true

negatives in the result thereby increasing its efficiency.

[2] "Combining lexicon and learning based approaches for concept level sentiment analysis": The approach in the

sentiment analysis is called pSenti. It uses both ML and Tokenisation approaches. It starts with Tokenisation and then carefully applies ML techniques wherever necessary. The paper provides an insight on the scope of the project beyond movie review classifications

[3] "Thumbs up: Sentiment classification using Machine

Learning Techniques."

The paper goes for contrasting the new procedures and the exemplary ones and after that amalgamates them into a learning system to comprehend the improvement there might be in the exactness. The investigation has been explicit to the issue of back to earlier extremity evaluation utilizing both relapse and grouping issue. Two distinct renditions of SentiWordNet have been utilized in this methodology. The outcome indicated 30 unique methodologies demonstrating that not a great deal of consideration had been given to the issues emerging in Sentiment examination and accentuates on the requirement for further examination. This paper has played a crucial role in highlighting the potential limitation there can be in the current system(s).

[4] "Review Classification using SentiWordNet":

The strategy proposed in the paper works in two stages: Interpretation stage and Polarity stage. In the translation stage for each word positive and negative scores are doled out dependent on the estimations of the normal for the passages. When this is finished with the second stage starts, the 'Sentiment Polarity Phase' Here the magnitude value of positive or negative scores are taken into consideration and the polarity is determined. The outcomes demonstrated an expansion of practically 10% in the exactness, where the connected term tallying strategy gave a precision of 56.77% this technique expanded it to 67.00%. The paper proposes the use of SentiWordNet which is used in this project for classification purpose.

## II. METHODOLOGY

The proposed system aims at expanding the proficiency of current assessment examination models by expanding the exactness of the outcomes produced. To achieve the goal it uses the above mentioned hybrid approach in a steady progression through the pipeline architecture.

The following are the various segments of the pipeline which are executed one after the other:

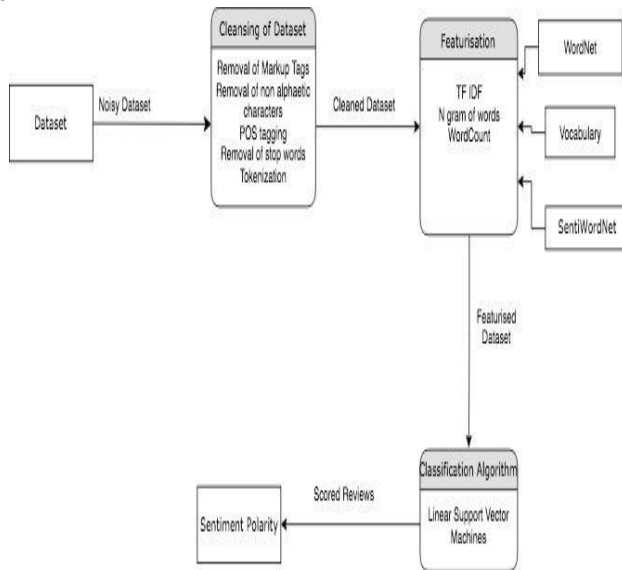
**Revised Manuscript Received on September 22, 2019.**

Devika Sharma, SRM Institute of science and Technology, e-mail: devikasharma\_ra@srmuniv.edu.in

Rajesh Babu, SRM Institute of science and Technology, e-mail: rajeshbabu.c@ktr.srmuniv.ac.in

### I. Obtaining the Dataset:

The dataset used in this paper is the IMDB dataset for movie reviews [5]. In the whole accumulation, close to 30



**Fig. 1.** Data Flow Diagram of the Methodology

surveys are taken into consideration any given film since audits for a similar motion picture will, in general, have associated appraisals. Further, the train and test sets contain a disjoint arrangement of movies, so no bias is introduced by retaining movie unique terms and their relatedness with labels. The dataset has already been divided evenly into training and testing sets. Every review is associated with a unique ID, so the dataset contains a unique ID followed by the review.

### III. CLEANING UP OF DATASET

The dataset in its raw form contains a lot of excess information that maybe a result of repetition or might just be unnecessary and irrelevant. This step thoroughly cleans the dataset by using various pre-processing methods such as removal of Markup tags, Removal of non-alphabetic characters, Tokenization, Removal of stop words, Parts of Speech, Removal of whitespaces and others.

### IV. FEATURISATION

A Feature is utilized principally for lessen vocabulary measure, clamor includes and expanded order exactness. Heterogeneous highlights made utilizing a blend of lexons like SentiWordNet, WordNet, and so on. . . and AI like sack of words, TF-IDF etc.” In this paper, we have made an exceptional feature utilizing

‘SentiWordNet, SentiWordNet Average on Review’: for each audit record, SentiWordNet scores dictated by ascertaining the normal of the aggregate of positive and

### V. CLASSIFICATION ALGORITHM

Sentiment Identification includes ordering surveys in content into classes dependent on opinion extremity. Extremity

in the survey content methods the condition of having two inverse audits. This extremity can be certain or negative. Here, in proposed approach record level opinion examination characterization utilized and here connected heterogeneous highlights to administered AI classifiers to learn and arrange content surveys into a positive and a negative classification. Here for this framework choosed an administered learning calculation since we can without much of a stretch a dataset large enough.

### VI. SENTIMENT POLARITY:

From regulated classifier, we can anticipate the feeling class name, it will be sure or negative class implies given audit content record will positive or negative.

### VII. PERFORMANCE

I. Choosing the right value of Inverse Regularization Parameter:

Inverse regularisation Parameter is a control variable that holds quality change of Regularization by being contrarily situated to the Lambda controller. It’s a term intended to control against Overfitting. We selected the value of C which gave us the highest accuracy. The accuracy score was calculated using the scikit learn library in Python (Fabian Pedregosa) using the following equation:

$$\text{accuracy}(y, \hat{y}) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}(\hat{y}_i = y_i)$$

The corresponding accuracy score for various values of C were calculated and based on the highest accuracy the value of C was selected.

#### A. Logistic Regression

Value of C	Accuracy for C
0.01	0.8968
0.05	0.89152
0.25	0.88976
0.5	0.88976
1	0.88928
Final Accuracy	0.89708

**Table 1:** Accuracy for various values of C

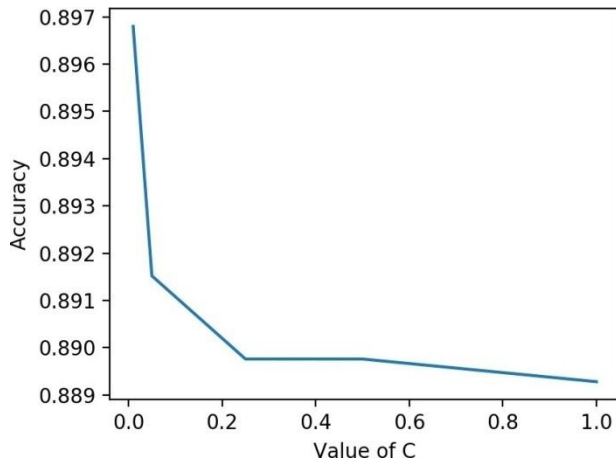


Fig 2: Accuracy graph for Logistic Regression

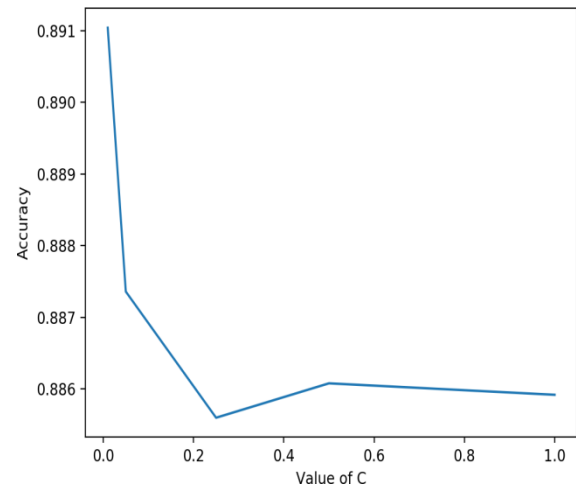


Fig 3 Accuracy graph for Naïve Baye's Classifier

#### B. Support Vector Machines:

Value of C	Accuracy for C
0.001	0.88784
0.005	0.89456
0.01	0.90064
0.05	0.89264
0.11	0.8928
Final Accuracy	0.90064

Table 2: Accuracy for various values of C

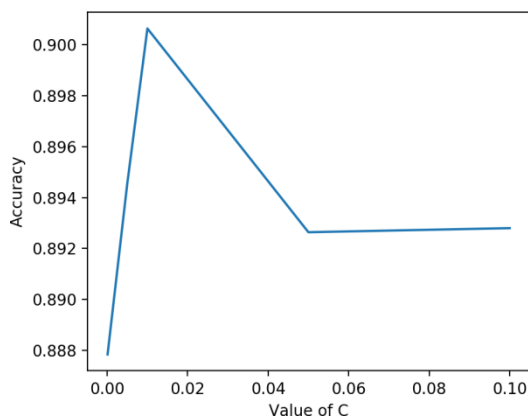


Fig 3: Accuracy graph for Support Vector Machines

#### C. Naïve Baye's Classifier:

Value of C	Accuracy for C
0.01	0.89104
0.05	0.88736
0.25	0.8856
0.5	0.88608
1	0.88592
Final Accuracy	0.89104

Table 3: Accuracy for various values of C

	Support Vector Machine	Naïve Baye's	Loistic Regression
Final Accuracy	0.90064	0.89708	0.89104

Table 4: Accuracy for various classifiers

## VIII. RESULT AND CONCLUSION

The paper presents a model that uses hybrid technique to identify the sentiment polarity of a given review. It not only calculates if the review is negative or positive, it also identifies the extent of negativity or positivity of the review y associating a score with the reviews. The model achieved an accuracy of 90.064% on identifying the sentiment polarity

## ACKNOWLEDGEMENT

I offer my unassuming thanks to Dr. Sandeep Sancheti, Vice Chancellor, SRM Institute of Science and Technology, for the offices stretched out for the project and his continous help.

I wish to express gratitude toward Dr. B. Amutha, Professor and Head, Department of Computer Science and Engineering, SRM Institute of Science and Technology, for her significant proposals and consolation all through the time of the task work.

I am very appreciative to our Academic Advisor Dr. K. Annapurani, Associate Professor, Department of Computer Science and Engineering, SRM Institute of Science and Technology, for her incredible help at all the phases of undertaking work. My inexpressible respect and thanks to my guide, Mr. C. Rajesh Babu, Department of Computer Science and Engineering, SRM Institute of Science and Technology, for providing me an opportunity to pursue my project under his mentorship.

I genuinely thank staff and understudies of the Computer Science and Engineering Department, SRM Institute of Science and Technology, for their assistance amid my exploration. At long last, we might want to thank my folks, my relatives and my companions for their unlimited love, consistent help and consolation

## REFERENCES

1. Archana Bandana (2018). Sentiment Analysis of Movie Reviews using Heterogeneous features
2. Andrius, Dell Zhang, and Mark Levene. "Combining lexicon and learning based approaches for concept-level sentiment analysis." Sivasankari,Bhanu Shri, Y.Bevish Jinila.
3. Alaa, and Mohamed Rohaim. "Reviews classification using sentiwordnet lexicon."
4. Pang Bo, Lillian Lee and Shivakumar Vaithyanathan. "Thumbs up? Sentiment classification techniques using Machine Learning Techniques"
5. Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts, "Learning word vectors for sentiment analysis"
6. Potts, Christopher. 2011. "On the negativity of negation. In Nan Li and David Lutz, eds., Proceedings of Semantics and Linguistic Theory" 20, 636-659.
7. H. M. Keerthi Kumar , B. S. Harish , H. K. Darshan "Sentiment Analysis on IMDb Movie Reviews Using Hybrid Feature Extraction Method"
8. M.Govindarajan "Sentiment Analysis of Movie Reviews using Hybrid Method of Naive Bayes and Genetic Algorithm"
9. Archana Bandana (2018). "Sentiment Analysis of Movie Reviews using Heterogeneous features "
10. Andrius, Dell Zhang, and Mark Levene. "Combining lexicon and learning based approaches for concept-level sentiment analysis." Sivasankari,Bhanu Shri, Y.Bevish Jinila.
11. Alaa, and Mohamed Rohaim. "Reviews classification using sentiwordnet lexicon."
12. Pang Bo, Lillian Lee and Shivakumar Vaithyanathan. "Thumbs up? Sentiment classification techniques using Machine Learning Techniques"
13. Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts, "Learning word vectors for sentiment analysis" Semantics and Linguistic Theory" 20, 636-659.
14. H. M. Keerthi Kumar , B. S. Harish , H. K. Darshan "Sentiment Analysis on IMDb Movie Reviews Using Hybrid Feature Extraction Method"
15. M.Govindarajan "Sentiment Analysis of Movie Reviews using Hybrid Method of Naive Bayes and Genetic Algorithm"
16. <https://www.imdb.com/>
17. Isa Maks,Piek Vossen " A lexicon model for deep sentiment analysis and opinion mining applications"
18. Theresa Wilson, Janyce Wiebe, Paul Hoffmann "Recognizing Contextual Polarity in Phrase-Level Sentiment Analysis"
19. B. Amutha "Performance Analysis of Firewall Based on SDN and OpenFlow"
20. B. Amutha "Study of Service Chain Optimization in Cloud Environment"