

SENTIMENT ANALYSIS OF MOVIE REVIEWS

Malini R

4th SEM, CS&E, AIT, Chikkamagaluru

Dr.Sunitha M.R

Professor, CS&E, AIT, Chikkamagaluru

Abstract: Sentimental analysis is one of the sub parts of opinion mining; it is the one of the new concepts of data mining. The online communication data consist of feedback in comments and reviews of particular topic that are posted on internet by internet users, where the analysis is focused on the extraction of emotions as a specific view or judgment on certain topic. Sentimental analysis system classifies text data into their respectively sentiments of positive polarity, negative polarity or neutral. In this paper, classification task of sentimental analysis of movie database is done. By using support vector machine algorithm the best accuracy is obtained.

Index terms: Sentimental analysis, polarity, sentiments, movie reviews, classifier.

I. INTRODUCTION

Sentimental analysis is rapidly increasing research area in the field of text mining. Posting online reviews on the different web sites has become an increasing popular way for people to share their opinions about specific product or services with other users. Sentimental analysis is the computational study of people judgment, attitudes and emotions towards an entity [1]. The entity can be represent individuals, events or certain topics. Opinion mining extracts and analyses people's opinion about an entity, while sentiment analysis finds the sentiments words expressed in a word of text document and it starts analysing it. Therefore the main goal of sentiment analysis is to find opinions, identify the sentiments they express, and classify their polarity. Sentimental analysis helps to find words that indicate sentiments and helps to understand the relationship between textual reviews and the significance of those reviews. One such domain of the reviews is the domain of movie which affects everyone from audience, film directors to the production company [1]. The movie reviews present on various websites are not formal reviews but they are rather very informal reviews and are unstructured form of grammar.

Sentiment analysis of twitter messages has picked up significance attention in the course of recent years. With the assistance of opinion mining, poor contents can be separated from high quality contents. It is possible that by utilizing available technologies, whether a movie has good opinion then bad opinion can be known and this helps the audience in their basic decision making. In this paper, reviews from the website are collected. The lexical approach is used for finding the overall polarity of the

movie reviews [5]. On the basis of entered reviews by the user it produces the results according to highest sentiments extracted from class of positive, negative and neutral. Finally the weak positive and strong positive features get converted into positive polarity and weak negative and strong negative features are converted into negative polarity. The result is graphically displayed in the form of positive and negative polarity.

II. LITERATURE SURVEY

A Large number of works have been carried out previously on opinion mining and sentiment analysis. Nagamma P et al. [1] proposed various data mining techniques for classification of movie audits and it likewise predicts the box office collection for the movie. The online movie review data collected from IMDB dataset [2], the box office collection and the success or failure of the movie is predicted based on the reviews. Pang et al. [3] applied the machine learning technique for classification of reviews present on IMDB movie reviews database, by forming the list of 14 keywords which are useful in finding the baseline for classification accuracy. The machine learning techniques like Naïve bayes, SVM, achieves higher accuracy over the baseline.

J. Erman et al. [4] studied three types of clustering algorithms namely K-Means, DBSCAN and Auto Class algorithm for the classification of network traffic problem. This study depends on the capacity of every algorithm for forming clusters having higher predictive power of a single traffic class and for deciding the capacity of every algorithm to produce small number of clusters that has numerous associations. The Auto Class algorithm is

compared with DBSCAN and K-Means algorithm and the result indicates that both K-Means and DBSCAN work faster than Auto Class algorithm. Turney et al. [5] studied the unsupervised learning algorithm for sentiment classification process. They determined the similarity of words with help of NEAR operation and developed a classifier for finding polarity result.

III. PROPOSED WORK

This section gives the description of the steps followed for the movie dataset mining for sentiment analysis. The movie review dataset is used and label to the polarity is provided as follows: positive, negative, positive-positive, positive-negative, negative-negative, and negative-positive.

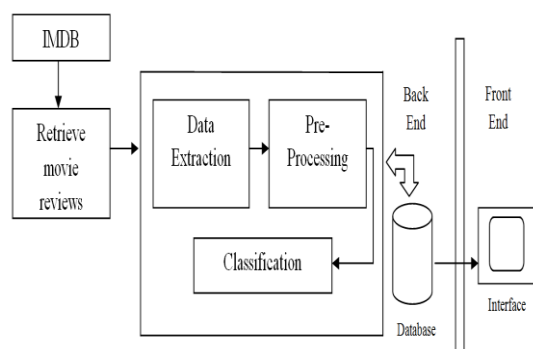


Fig 1. Architecture diagram for sentiment analysis

Fig.1 demonstrates the architecture diagram for sentiment analysis. The movie reviews are collected from IMDB database to determine the sentiment of each review. Support vector machine strategy is utilized to analyze the sentiment of the tweets. The reviews are classified as positive review, negative review or neutral review.

A. Input Data

The input data is in the form of reviews from the movie review dataset. Particular movie is selected from the dataset and reviews regarding that movie are displayed on web page. After releasing of any new movie the reviews of that movie are added to the dataset.

B. Processing

The text pre-processing techniques are:

Tokenization: The data present in the text document contains block of characters called tokens. These content reports are isolated as tokens and utilized for further handling of information.

Removal of Stop Words: A web search tool or other natural language processing system may contain collection of stop-records, or it might contain a solitary stop-list. Most of the more frequently used stop words in English are “an”, “a”, “of”, “the”, “you”, “and” these are some words which do not carry any meaning. Hence, those words which appear too often that support no information for the task are removed.

C. Classification

Many approaches are mainly classified into two categories namely lexicon based approach and machine learning based approach. The lexicon based approach is used for finding the overall polarity of movie reviews. The well known classifier namely Support Vector Machine (SVM) technique is used for sentiment classification. The classification is done with Support Vector Machine (SVM) classifier to determine the sentiment labels for a machine and to predict the class of a movie reviews when never it arrives in the form of positive or negative polarity. Classification Rate or Accuracy is given by the equation 1 as shown below:

$$\text{Accuracy} = \frac{TP+TN}{TP+TN+FP+FN} \quad (1)$$

where TP is True Positive, TN is True Negative, FP is False Positive, FN is False Negative.

IV. Results

In this work, the classification results are obtained by applying support vector machine algorithm to determine the polarity of the reviews.

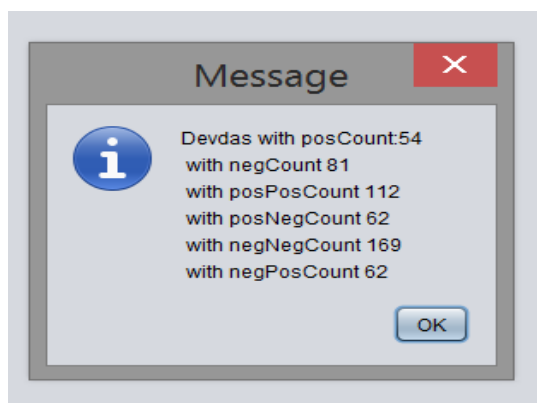


Fig 2: Classification count for the movie

Fig.2 demonstrates the sentiment classification count of specific movie utilizing Support Vector Machine classifier. A supervised machine learning algorithm that can be used for classification is the Support Vector Machine. The reviews are arranged for a specific movie by Support Vector Machine classification. Positive, negative, positive-positive, positive-negative, negative-negative, and negative-positive are the sentiment classification count determined for every movie.

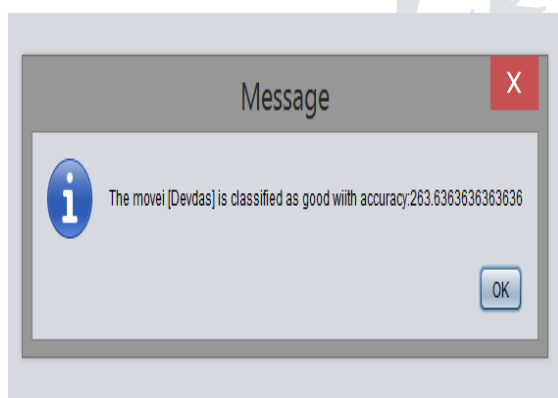


Fig 3: Support vector machine classification accuracy

Fig.3 demonstrates the classification accuracy of support vector machine algorithm. Support-vector machine creates a hyper plane in an infinite-dimensional space, which can be used for classification, regression purposes. The reviews are classified for a specific movie with the help of support vector machine classification. The accuracy of classification algorithm is calculated for each movie.

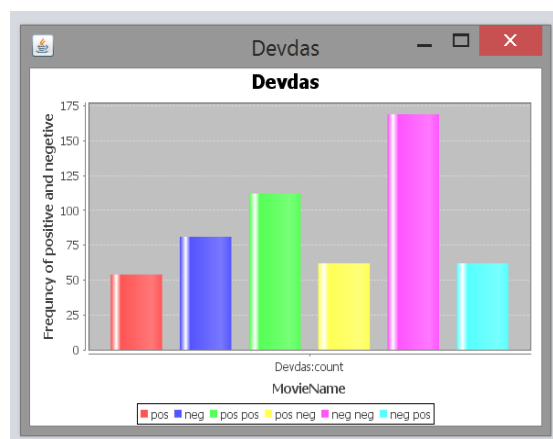


Fig 4: Graph for sentiment analysis of a particular movie

Fig.4 shows the bar chart for sentiment analysis of particular movie using Support Vector Machine algorithm. The sentiment analysis graph for a particular movie is plotted with x-axis denoting the movie name and y-axis denoting the frequency of positive and negative counts for the movie. The red bar indicates the positive count, blue bar indicates the negative count, green bar indicates the positive-positive count, yellow bar indicates the positive-negative count, pink bar indicates the negative-negative count and light blue bar indicates the negative-positive count for the movie.

V. Conclusion

Analysis of sentiment of twitter is crucial for the sake; it is very tough to analyze sentimental argument pattern tweets. And also it is hard as the repeated characters, slang words, whitespaces, misspellings are present. In order to supervise these problems, Natural Language Processing is used. Pre-processing is performed on individual tweet, before NLP is adopted. First phase is to find misspellings, repeated characters, and special characters. The text containing unwanted data are removed. Instantly, extricated movie text reviews are converted into ordinary text. Later, movie reviews extricated, which is in ordinary text pattern free from hash tags is created.

In order to perform tweets classifications, there are several machine learning classifiers. Support vector machine executes well and also provide higher accuracy as shown in the results. So, as training data is increased, the classification accuracy can also be increased. The outcome of this project is that Natural Language Processing operates superior for movie reviews related tweets.

ACKNOWLEDGMENT

This is a part of post graduate project work and I represent my sincere gratitude to all my teachers for their constant guidance throughout the work and providing excellent atmosphere for dissertation work. In future I would like to evaluate the effectiveness of the proposed sentimental classification features and techniques for other tasks such as sentiment classification. I would like to apply in depth concepts of SVM for better prediction of the polarity of the document. We would like to extend this technique on other domains of opinion mining likes newspaper articles, product reviews, political discussion forums etc.

REFERENCES

- [1] P.Nagamma, Pruthvi H.R, Nisha K.K, Carlos Soares, "An Improved Sentiment Analysis of Online Movie Reviews", IEEE 2015, International conference on Computer and Information Technology.
- [2] IMDB Dataset: www.imdb.com/india/top-rated-indian-movies/
- [3] B. Pang, L. Lee, and S. Vaithyanathan, "Thumbs up? Sentiment classification using machine learning techniques," in Proceedings of the ACL-02 conference on Empirical methods in natural language processing-Volume 10. Association for Computational Linguistics, 2002, pp. 79–86.
- [4] J. Erman, M. Arlitt, and A. Mahanti, "Traffic classification using clustering algorithms," in Proceedings of the 2006 SIGCOMM workshop on Mining network data. ACM, 2006, pp. 281–286.
- [5] Turney, Peter, and Michael L. Littman, "Unsupervised learning of semantic orientation from a hundred-billion-word corpus", (2002).
- [6] Baccianella, Stefano, Andrea Esuli, and Fabrizio Sebastiani, "SentiWordNet 3.0: An Enhanced Lexical Resource for Sentiment Analysis and Opinion Mining", LREC.Vol. 10. 2010.
- [7] Prabowo, Rudy, and Mike Thelwall, "Sentiment analysis: A combined approach", Journal of Informetrics 3.2 (2009): 143-157.
- [8] Rui Yao and J. Chen "Predicting movie sales revenue using online reviews", In GeC, 2013, pp. 396-401.
- [9] Ion Smeureanu, Cristian Bucur, "Applying Supervised Opinion Mining Techniques On Online User Reviews", Informatica Economică, 2012.
- [10] Singh, V. K., et al. "Sentiment analysis of movie reviews: A new feature based heuristic for aspect-level sentiment classification", Automation, and the Computing, Communication, Control and Compressed Sensing (iMac4s), International Multi-Conference on. IEEE, 2013.