

AN IN-DEPTH STUDY OF THE SENTIMENT OF IMDB (INTERNET MOVIE DATABASE) MOVIE AND RATINGS

Kanishka Kashyap

Vandana International Sr. Sec. School, New Delhi

ABSTRACT

An area of artificial intelligence known as sentiment analysis focuses on interpreting data to convey human feelings and opinions. This paper focuses on the IMDB movie review sentiment analysis, and Stanford University provided the dataset. We Analyse how the viewer articulates their positive or negative opinion. Because, at this time, we don't have an accurate system that can frame the structures in the study of random slang movies, the N-gram method was used. Along with that, we make a selection of standard features and use them for training multiple-label classifiers to tag reviews of movies accurately and further, choosing the most suitable classifier for our domain query by various categories approach comparison. Our process, which makes use of separation techniques, is 83% accurate.

INTRODUCTION

"What other people think" has always been the most secure method for making well-informed decisions regarding voting, service seeking, and shopping. The Internet is now accessible to more people than just friends, family, and neighbours; It now includes newcomers to our lives. A feature known as sentiment tells us whether an opinion is neutral, positive, or negative about a specific subject. Alternatively, a lot of people express their reviews on the other website. [1] "Logan is a nice movie, highly recommended, 10/10," for instance, expressing admiration for the movie Logan and its title. For instance, in the text "I'm stunned that such countless individuals put Logan in their #1 motion pictures, I felt a decent watch however not unreasonably great," the creator's sentiments about the film may be great, yet not a similar message. The subject of deep examination was an item or administration freely open on the web surveys.

This paper aims to classify movie reviews [3]. This may account for the frequent use of the terms "theoretical analysis" and "spiritual analysis" together; However, we think it is more accurate to think of emotions as ideas with strong emotional attachments. Because we used sentiment words a lot in the study, so we devised a way to describe the movie's polarity using these words.

LITERATURE REVIEW

A. Natural Language Processing (NLP)

NLP is the ability of a computer program to recognize the natural form of human speech. NLP is included in artificial intelligence. NLP applications are notoriously difficult to develop because computers typically require people to "speak" to them using a limited number of

written voice commands or the language of an accurate, well-structured system. On the other hand, human speech is only precise sometimes.

B. Sentiment Analysis

Also known as opinion mining, it is a type of NLP that examines composed correspondence to decide its general state of mind. In business, it's common practice to identify and categorize ideas about a product, service, or idea. Machine learning, artificial intelligence and Data analysis (AI) are necessary for opinion-mining text and sensitive data.

C. Pre-Processing Text and Normalization

Cleaning, pre-processing, and normalizing the text in order to bring text artifacts like words and phrases into some degree of format is one of the most important steps before participating in the cutting-edge engineering and modelling process.

Positioning over the whole text is made conceivable by this, which makes significant usefulness and decreases the commotion that can be imported because of inert characters, exceptional element characters, XML and HTML labels, and different variables. [9].

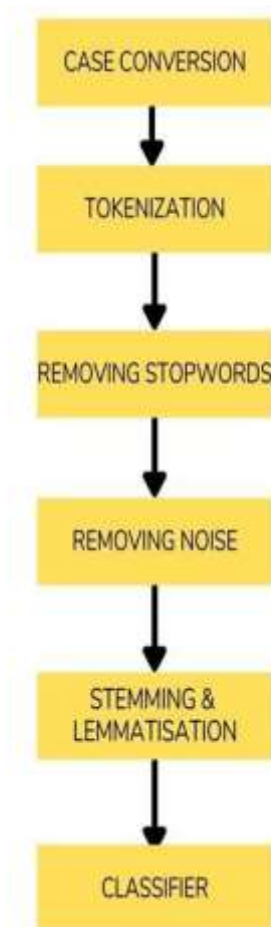


Fig. 1 Steps for Text Pre-processing

The most important parts of our text normalization pipeline are as follows:

- 1) Text for Cleaning: When evaluating our emotions, unnecessary content like HTML tags frequently appears in our text and has little significance. As a result, before the functions are removed, we must ensure they are disabled. The BeautifulSoup library excellently provides the necessary facilities.
- 2) Eliminating Highlighted Characters: In our database, we work with English language changes; however, we must ensure that computerized characters and other formats are translated into ASCII characters.
- 3) Contractions that get bigger: In English, problems are shortened forms of words or letter clusters. These shortened versions are made by omitting particular letters and sounds from actual words or phrases. English usually takes the place of vowels.
- 4) Removing unique characters: Removing special characters and icons, which frequently impart an additional sound to meaningless text, is yet another crucial aspect of text processing and design.
- 5) Lemmatization and stemming: Word titles are always the simplest form of terms created by adding prefixes and suffixes to the root of existing words.
- 6) Stop Words Removal: Stop words are those words which is of less use or of no use, like is, am, are the. These words are important to complete the sentence.

D. Methods for Feature Weighting

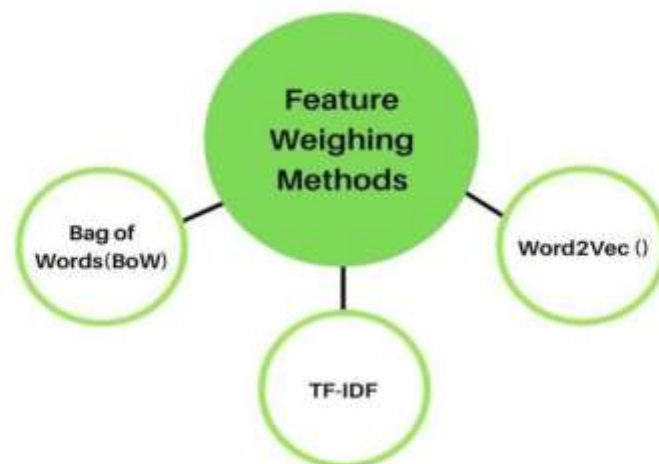


Fig. 2 Feature Weighing Methods

The selection of text-related functions based on the number of continuous values is known as feature weighting. We tried to test the viability of weight reduction approaches in the feeling of profound segregation. Additionally, we attempted to determine the impact of the split

algorithm. From among the exceptions, we selected the following three weight-measuring methods:

1) The Bag of Words: BoW is used to calculate the word that appears in a text again and again; the paper is compared to this, and comparisons are made using various applications. It gets the vector VT following the text T , where VT_i shows the times, and the word shows up in the text. T . D . is a lexicon that has been revised to include all words.

2) TF-IDF: This indicates a Document Frequency that can be modified repeatedly if it demonstrates the meaning of the expression in the context.

When evaluating diversity mining applications, this is taken into account. Please increase the number of times a word is used in the text to eliminate its use in the corpse.

3) W2V, or Word2Vec: Word2Vec is used to develop term embedding. Word meaning is formed by this model using two-layer neural networks. It produces a space vector of words from a large text corpus as its input. A spatial vector is assigned to each name in the corpus. Names are shared because the corpus's general meaning is similar to the vector space. Sentences are created from the definitions before running the Word2Vec algorithm.

E. Text Classifiers

In sensory research, text separators are used to determine which movie reviews receive the most votes. Consequently, we employ the methods outlined below to investigate the impact of dividers on updates.

1) RF Random Forest: It's a way for individuals to meet. The approach used to improve performance is solved with differentiated integration approaches. A group of good students has been formed from several vulnerable students. A guarded machine first reads the decision tree. [6]

2) NB, or Naive Bayes: utilizing the Bayesian method. The presence of another element is unaffected by the appearance of one in the classroom. The NB model is easy to use and works well with a lot of data. The back odds are calculated using the Bayes theorem.

3) The posterior probability of the class-given predictor is calculated as $P(c/x) = (P(c/x)*P(c))/P(x)$.

$P(c)$ is the class's prior probability.

Probability = $P(x/c)$.

$P(x)$ is the predictor's prior probability.

4) Support Vector Machine (SVM): It is an algorithm for supervised machine learning. It is frequently utilized for the withdrawing side as well as partitioning. The idea is to divide the database into two sections to locate a hyperplane. The space between the data of interest and the hyperplane is known as the limit. The objective is to determine the absolute limit of the

potential for well-distributed new data. The accuracy of smaller collections is higher, but large-scale data sharing takes longer.

5) Neural Network (NN): They simultaneously work on an individual record. Learned by comparing record divisions of well-known classifications. It isn't easy to train and tune a large database. Mistakes from the main emphasis are taken care of in the nonstop stream, lessening blunders. Back delivery, which alters the link's weight through forward simulation, is frequently utilized for training networks that have demonstrated high performance. It enhances the precision of various data.

6) KNN, or K-Nearest Neighbor, uses every possible scenario for segmentation. New classifications are established using the Euclidean parallel scale. Pattern analysis and statistical equations both make use of it. If the intricacy of the training database expands, KNN turns out to be more perplexing.

7) SGD, or stochastic gradient descent: Another name for it is the rising gradient. Adjust parameter values by human training values. This technique is helpful when an improvement calculation is searching for boundaries. Rather than re-establishing coefficients after a progression of boundaries, it is finished for each preparing occasion. [7] Because the gradient drop algorithm anticipates each data state, processing large data sets takes longer. In short, massive data are positively affected by the reduction of stochastic gradients.

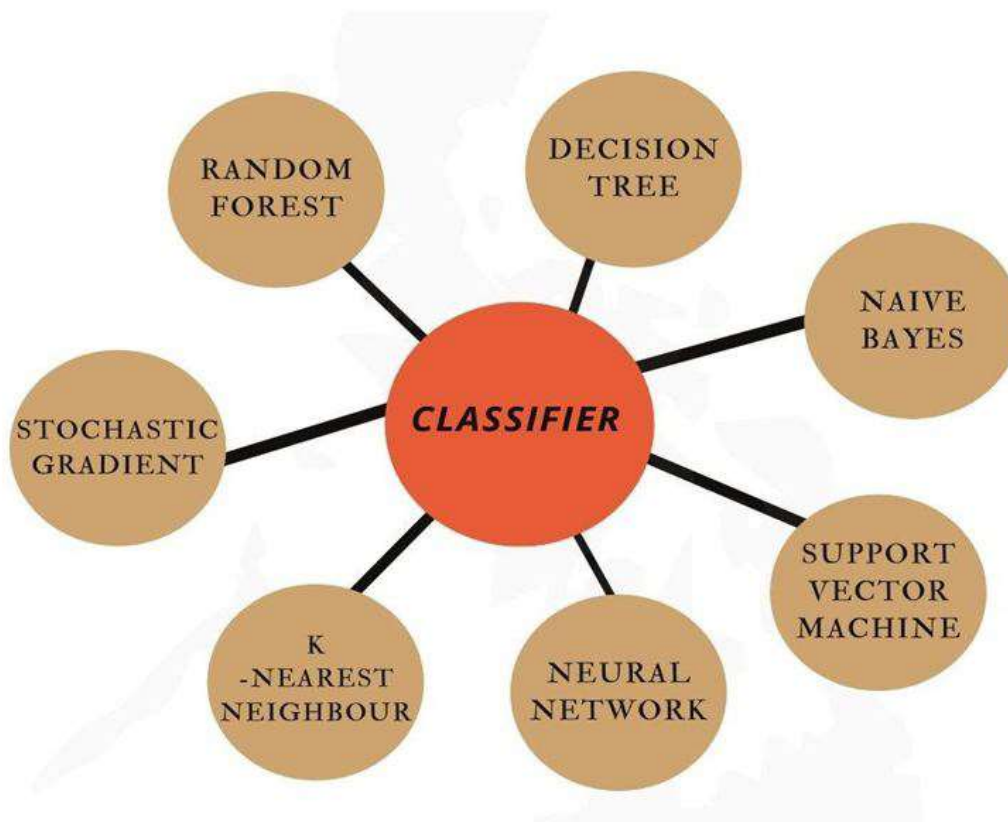


Fig 3 Classifiers

RESULTS AND DISCUSSION

We used precision, recall, accuracy, and F1 ranking as performance metrics [17].

A. Accuracy It is a distinct subset of the total number of emotions. It is used to determine a machine's effectiveness.

Precision is a test of how many feelings have been correctly judged as right as opposed to the overall amount of well-separated emotions (P) Precision = $TP/(TP+FP)$ C. Recall is the overall number of sensors in the absolute total that are correctly predicted (R) Recall = $TP/(TP+FN)$

$$F1=(2*P*R)/(P+R)$$

1) The recall of the SGD classifier was 86.7%, whereas the accuracy, precision, and F1 ratings of the BoW + SVM combination were 84%, 83.33%, and 85%, respectively.

2) We achieved good accuracy and precision for the SGD classifier by combining TFIDF with other classifiers (75.1% and 78.1%, respectively). The recall and F1 scores of the TF-IDF + SVM combinations were 75.3% and 75.3%, respectively.

3) When Word2Vec was paired with other classifiers, the RF classifier had a high precision of 80.4%, whereas the Word2Vec + SGD combination had an accuracy of 82.2%, a recall of 89.0% and an F1 score of 83.6%.

4) We found the following when comparing precision:

- a) as far as BoW, SVM is an area of strength for extremely (per cent)
- b) In TF-IDF, SGD has a high rate (80.9%).
- c) Random Forest receives a high score in Word2Vec (76.6%).

CONSTRAINTS

The primary obstacles to sentiment analysis are as follows: -

- 1) Identify an entity with a name: What topic is the user discussing? For instance, does the expression "300 Spartans" refer to a movie or a group of ancient Greeks?
- 2) A Solution to Anaphora: Anaphora resolution refers to the difficulty of determining what a pronoun or noun phrase refers to. What precisely does "It" imply? We ate dinner and saw the movie together; It was awful."
- 3) Sorting: Which person or thing is the subject and object of this sentence, as well as the noun and adjective?
- 4) Arrogance: If you don't know the author, you won't know whether "bad" means "bad" or "good."

5) Twitter: On Twitter, there are issues with capitalization, abbreviations, improper pronunciation, punctuation, and syntax.

CONCLUSION

The topic of sentiment analysis is the focus of this paper. There are three factors. To identify emotions in the IMDB database of movie reviews, we integrate the Bag of Words, Term Frequency-Inverse Document Frequency, and Word2Vec analysis with various classifiers, such as the Decision tree, Random Forest, Naive Bayes classifier, Support vector machine, and checks. Word2Vec with SGD effectively addresses the IMDB database segmentation issue, as demonstrated by the test results. This form was the best of all combinations because it had the highest accuracy score of 89 per cent and the highest F1 score of 83.6%. The analysis of online feedback from social media platforms and sophisticated algorithms could be part of this study.

REFERENCES

- [1] NLTK: The Natural Language Toolkit. Edward Loper and Steven Bird Department of Computer and Information Science University of Pennsylvania, Philadelphia, PA 19104-6389, USA
- [2] Steven Bird NLTK Documentation Release 3.2.5. (English) Sep 28, 2017
- [3] Edward Loper, Steven Bird, NLTK: the natural language toolkit, Proceedings of the ACL 2004 on Interactive poster and demonstration sessions, p.31-es, July 21-26, 2004, Barcelona, Spain
- [4] Naive Bayes Scikit-Learn.org. (English) Sep 28, 2017
- [5] Cortes, Corinna; Vapnik, Vladimir N. (1995). "Support-vector networks". *Machine Learning*. 20 (3): 273–297. CiteSeerX 10.1.1.15.9362
- [6] Guerini Marco Lorenzo Gatti Marco Turchi Sentiment analysis: How to derive prior polarities from SentiWordN et vol. 5843 2013.
- [7] Bhoir Purata Shilpa Kolte "Sentiment analysis of movie reviews using lexicon approach" 2015 IEEE International Conference on Computational Intelligence and Computing Research (ICIC) pp. 1-6 2015.