

Sentiment Analysis on Movie Reviews

S. Prathap¹, Sk. Moinuddin Ahmad²

CSE Department, Anil Neerukonda Institute of Technology and Sciences, Visakhapatnam, Andhra Pradesh,
India

CSE Department, Gudlavalleru Engineering College, Gudivada, Andhra Pradesh, India

ABSTRACT

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. We try to focus our task of sentiment analysis on IMDB movie review database. Sentiment Analysis is a process of extracting information from large amount of data, and classifies them into different classes called sentiments. Python is simple yet powerful, high-level, interpreted and dynamic programming language, which is well known for its functionality of processing natural language data by using NLTK (Natural Language Toolkit). NLTK is a library of python, which provides a base for building programs and classification of data. NLTK also provide graphical demonstration for representing various results or trends and it also provide sample data to train and test various classifier respectively. Sentiment classification aims to automatically predict sentiment polarity of users publishing sentiment data. Traditional classification algorithm can be used to train sentiment classifiers from manually labeled text data. We directly apply a classifier trained to the domain to the performance will be very low due to the difference between these domains. In this work, we develop a general solution to sentiment classification when we do not have any labels in target domain but have some labeled data in a different domain, regarded as source domain. In this project, we attempt not to only detect sarcasm in text and made pilot Model Sarcasm with Naive Bayes using TFIDF feature vectors.

Keywords : NLTK, IMDB, NLTK, NLTK, Python Programming

I. INTRODUCTION

Text can be categorized in two types based on its properties in terms of text mining : 'subjectivity' and 'polarity'. The focus of our project is to find the polarity of the text which means that we are interested in finding if the sentence is positive or negative. We use machine learning techniques classify such sentences and try to find answers to the following questions:

1. Machine learning techniques their purpose, which one out of them performs the best and which techniques are better than the others.
2. Advantages and disadvantages of traditional machine learning techniques for sentiment analysis.
3. The difficulty in the task of extracting sentiment from short comments or sentences can be as compared to the traditional topic based text classification.

Sentimental analysis is a hot topic of research. Use of electronic media is increasing day by day. Time is money or even more valuable than money therefore instead of spending times in reading and figuring out the positivity or negativity of text we can use automated techniques for sentiment analysis. Sentiment analysis is used in opinion mining. The applications of sentiment analysis are broad and powerful. The ability to extract insights from social data is a practice that is being widely adopted by organisations across the world. Shifts in sentiment on social media have been shown to correlate with shifts in the stock market.

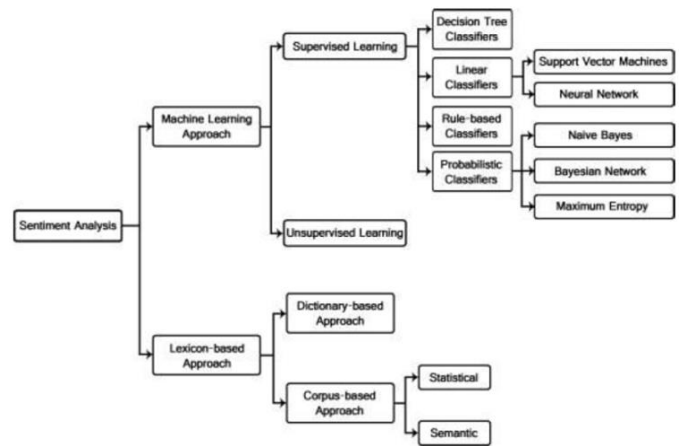
Sentiment analysis is becoming one of the most profound research areas for prediction and classification. Automated sentiment analysis of text is used in fields where products and services are reviewed by customers and critics. Thus, sentiment analysis becomes important for businesses to draw

a general opinion about their products and services. Our analysis helps concerned organizations to find opinions of people about movies from their reviews, if it is positive or negative. One can in turn formulate a public opinion about a movie. Our goal is to calculate the polarity of sentences that we extract from the text of reviews. We will model sentiment from movie reviews and try to find out how this sentiment matches with the success for these movies. In other words, if a movie review is positive, negative or neutral. But this task can be difficult and tricky. Consider a sentence “the movie interstellar was visually a treat but the story line was terrible”. Now one can clearly see how categorizing this sentence as negative, positive or neutral can be difficult. The phrases “visually a treat” and “story line was terrible” can be considered positive and negative respectively but the degree of their positiveness’ and ‘negativeness’ is somewhat ambiguous. We use a score for common positive and negative words and use this score to calculate the overall sentiment of a sentence.

Based upon the above reasons, this project applies the Naive Bayes classifier machine learning techniques to predict polarity of movie reviews as positive or negative by understanding meaning and relationship between the words. Mining the movie reviews and generating valuable metadata provides an opportunity to understand the general sentiment around movies in an independent way. The project is implemented using Python Programming Language and machine learning libraries of Python to predict sentiment of movie reviews as positive or negative using Naive Bayes machine learning algorithm.

II. METHODS AND MATERIAL

1. Classification



2. Models

There are many ways to implement Sentiment Analysis. Ultimately, it is a text classification problem and can be broken down into two main areas

Supervised models: This technique involves the construction of a “Classifier” and the problem has been studied intensively. The Classifier is responsible for categorizing texts into either a positive, negative or neutral polarity. The three main classification techniques are:

- i) Naive Bayes
- ii) Maximum Entropy
- iii) Support Vector Machines (SVM)

Unsupervised models: Unsupervised Learning has three steps.

1. Implement POS tagging (Part of Speech), then, two consecutive words are extracted to identify if their tags conform to given patterns.
2. Estimate the sentiment orientation (SO) of the extracted phrase.
3. Compute the average SO of all phrases that were extracted in terms of positive or negative.

Sentiment Analysis Methods

The sentiment classification approaches can be classified in

- (i) Machine Learning : The machine learning approach is used for predicting the polarity of sentiments based on trained as well as test data sets
- (ii) Lexicon based : Lexicon based approach does not need any prior training in order to mine the

data. It uses a predefined list of words, where each word is associated with a specific sentiment.

(iii) Hybrid approach : combination of both the machine learning and the lexicon based approaches has the potential to improve the sentiment classification performance.

3. Feature Extraction

Text Analysis is a main application field for mechanism learning processes. However the raw information, an order of symbols cannot be fed straight into the algorithms themselves as maximum of them expect arithmetical feature paths with a fixed size somewhat than the raw text forms with variable length. In imperative to address this, sickie-learn offers utilities for the most mutual ways to extract numerical structures out of texts, as follows:

- ✓ Tokenizing the strings and giving an integer id for each imaginable token, for example by using
- ✓ white-spaces & punctuation as symbolic separators.
- ✓ Counting the existences of tokens in each document.
- ✓ Regulating and weighting with diminishing importance tokens that occur in the majority of samples / forms.

4. Existing Sentiment Technnique:

Naive Bayes Algorithm: Bayes theorem named after Rev. Thomas Bayes. It works on conditional probability. Conditional probability is the probability that something will happen, *given that something else* has already occurred. Using the conditional probability, we can calculate the probability of an event using its prior knowledge. Naive Bayes model is easy to build and particularly useful for very large data sets Below is the formula for calculating the conditional probability

$$P(c|x) = \frac{P(x|c)P(c)}{P(x)}$$

Likelihood
Class Prior Probability
Posterior Probability
Predictor Prior Probability

$$P(c|X) = P(x_1|c) \times P(x_2|c) \times \dots \times P(x_n|c) \times P(c)$$

Where

- ✓ P(c) is the probability of hypothesis c being true. This is known as the prior probability.
- ✓ P(x) is the probability of the evidence (regardless of the hypothesis).
- ✓ P(x|c) is the probability of the evidence given that hypothesis is true.
- ✓ P(c|x) is the probability of the hypothesis given that the evidence is there.

Working:

Let's understand it using an example. Below I have a training data set of weather and corresponding target variable 'Play' (suggesting possibilities of playing). Now, we need to classify whether players will play or not based on weather condition. Let's follow the below steps to perform it.

Step 1: Convert the data set into a frequency table

Step 2: Create Likelihood table by finding the probabilities like Overcast probability = 0.29 and probability of playing is 0.64.

Step 3: Now, use Naive Bayesian equation to calculate the posterior probability for each class.

The class with the highest posterior probability is the outcome of prediction.

Disadvantages

The first disadvantage is that the Naive Bayes classifier makes a very strong assumption on the shape of your data distribution, i.e. any two features are independent given the output class. Due

to this, the result can be (potentially) very bad hence, a “naive” classifier

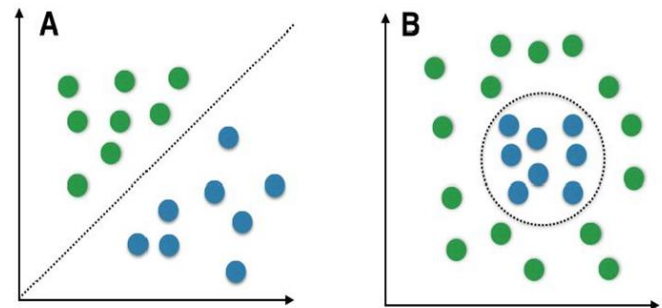
Another problem happens due to data scarcity. For any possible value of a feature, you need to estimate a likelihood value by a frequentist approach. This can result in probabilities going towards 0 or 1, which in turn leads to numerical instabilities and worse results

III. METHODS AND MATERIAL

Proposed System

Naive Bayes Classification: Naive Bayes classifiers are linear classifiers that are known for being simple yet very efficient. The probabilistic model of naive Bayes classifiers is based on Bayes’s theorem, and the adjective naive comes from the assumption that the features in a dataset are mutually independent. In practice, the independence assumption is often violated, but naive Bayes classifiers still tend to perform very well under this unrealistic assumption. Especially for small sample sizes, naive Bayes classifiers can outperform the more powerful alternatives. Being relatively robust, easy to implement, fast, and accurate, naive Bayes classifiers are used in many different fields. Some examples include the diagnosis of diseases and making decisions about treatment processes the classification of RNA sequences in taxonomic studies, and spam filtering in e-mail clients. However, strong violations of the independence assumptions and non-linear classification problems can lead to very poor performances of naive Bayes classifiers. We have to keep in mind that the type of data and the type of problem to be solved dictate which classification model we want to choose. In practice, it is always recommended to compare different classification models on the particular dataset and consider the prediction performances as well as computational efficiency. In the following sections, we will take a closer look at the probability model of the naive Bayes classifier and apply the concept to a simple toy

problem. Later, we will use a publicly available SMS (text message) collection to train a naive Bayes classifier in Python that allows us to classify unseen messages as spam or ham.



Linear (A) vs non linear Problems (B)

Random samples for two different classes are shown as colored spheres, and the dotted lines indicate the class boundaries that classifiers try to approximate by computing the decision boundaries. A non-linear problem (B) would be a case where linear classifiers, such as naive Bayes, would not be suitable since the classes are not linearly separable. In such a scenario, non-linear classifiers (e.g., instance-based nearest neighbor classifiers) should be preferred.

Class-conditional Probabilities

One assumption that Bayes classifiers make is that the samples are i.i.d. The abbreviation i.i.d. stands for “independent and identically distributed” and describes random variables that are independent from one another and are drawn from a similar probability distribution. Independence means that the probability of one observation does not affect the probability of another observation. One popular example of i.i.d. variables is the classic coin tossing: The first coin flip does not affect the outcome of a second coin flip and so forth. Given a fair coin, the probability of the coin landing on “heads” is always 0.5 no matter of how often the coin is flipped.

An additional assumption of naive Bayes classifiers is the conditional independence of features. Under this

naive assumption, the class-conditional probabilities or (likelihoods) of the samples can be directly estimated from the training data instead of evaluating all possibilities of x . Thus, given a d -dimensional feature vector x , the class conditional probability can be calculated as follows:

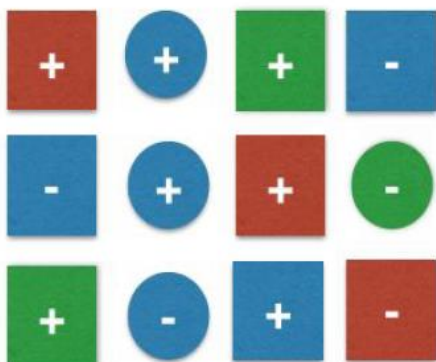
$$P(x|\omega_j) = P(x_1|\omega_j) \cdot P(x_2|\omega_j) \cdot \dots \cdot P(x_d|\omega_j) = \prod_{k=1}^d P(x_k|\omega_j)$$

Here, $P(x|\omega_j)$ simply means: "How likely is it to observe this particular pattern x given that it belongs to class ω_j ?" The "individual" likelihoods for every feature in the feature vector can be estimated via the maximum-likelihood estimate, which is simply a frequency in the case of categorical data:

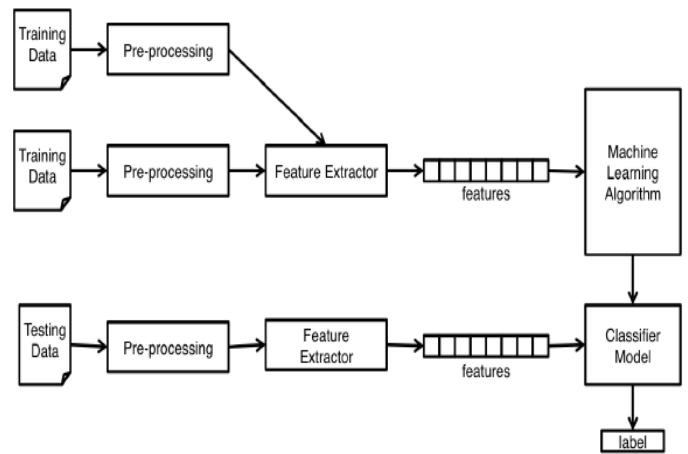
$$P^{(x_i | \omega_j)} = \frac{N_{x_i, \omega_j}}{N_{\omega_j}} \quad (i = 1, \dots, d)$$

- ✓ N_{x_i, ω_j} : Number of times feature x_i appears in samples from class ω_j .
- ✓ N_{ω_j} : Total count of all features in class ω_j .

Multinomial Naive Bayes - A Toy Example:
After covering the basics concepts of a naive Bayes classifier, the posterior probabilities and decision rules, let us walk through a simple toy example based on the training set shown in



Example for Multinomial NB , the toy example Architecture



IV.CONCLUSION

In this monitoring and management for green environment.

Cite this article as :

V. REFERENCES

WEB SITES:

- [1]. <http://pythonforengineers.com/build-a-sentiment-analysis-app-with-movie-reviews/>
- [2]. <https://www.programiz.com/python-programming/methods/built-in/map>
- [3]. <http://vpython.org/contents/docs/vector.html>
- [4]. <http://www.nltk.org/book/ch07.html>
- [5]. <http://www.nltk.org/book/ch01.html>
- [6]. <https://link.springer.com/article/10.1023/A:1007673816718>
- [7]. <https://en.wikipedia.org/wiki/Metaphone>
- [8]. <https://taku910.github.io/crfpp/>

BOOKS

- [9]. NLTK Book by Steven Bird textbook.
- [10]. Think Python by Allen.B.Downey 2nd edition.
- [11]. Machine Learning by Vinod Chandra 2 nd edition.