

Sentiment Analysis of College Reviews

¹Omkar Borade, ²Kaushik Gosavi, ³Ajay Shinde, ⁴Avinash Gowda

¹B.E. Student, ²B.E. Student, ³Professor

¹Name of Department of 1st Author,

¹Terna Engineering College, New-Mumbai, India

Abstract—The use of reviews has created many opportunities for people to publicly voice their opinions. Reviews in the internet could be in millions for services which make it difficult to track and understand customer opinions. Sentiment analysis is an emerging area of research to extract the subjective information to track and understand customer opinions. The reviews provide accessible and plentiful data for relatively easy analysis for a range of applications. This system seeks to apply and extend the current work in the field sentiment analysis on college reviews data retrieved from another websites using web mining technique. Naive Bayes and decision list classifiers are used to tag a given review as positive or negative. The features, such as bag-of- words and bigrams, are compared to one another in their effectiveness in correctly tagging reviews. Recent studies analyzed this reviews and found that it includes information useful for college, such as user requirements, ideas for improvements, user sentiments about specific features, and descriptions of experiences with these features.

Index Terms—Review, Sentiment Analysis, Naïve Bayes. (keywords)

I. INTRODUCTION

Reviews means the text which is given by user related with our system service. This review is useful for get detail about our system means what is better in our system and which changes is require for make our system better. Using these reviews we can take opinion of the user of our system. A user review refers to a review written by a user for a product or a service based on her experience as a user of the reviewed service or product. Most of college website gets reviews from users so the college system see view services is better and which services is require changes. Using this reviews college can make changes in services. Popular sources for consumer reviews are e-commerce sites like Amazon or Zappos, and social media sites like Trip Advisor and Yelp. E-commerce sites often have consumer reviews for products and sellers separately. Usually, consumer reviews are in the form of several lines of texts accompanied by a numerical rating. This text is meant to aid in shopping decision of a prospective buyer. A consumer review of a product usually comments on how well the product measures up to expectations based on the specifications provided by the manufacturer or seller. It talks about performance, reliability, quality defects, if any, and value for money.

Sentiment analysis can be defined as a process that automates mining of attitudes, opinions, views and emotions from text, speech, tweets and database sources through Natural Language Processing (NLP). Sentiment analysis involves classifying opinions in text into categories like “positive” or “negative” or “neutral”. It’s also referred as subjectivity analysis, opinion mining, and appraisal extraction. The words opinion, sentiment, view and belief are used interchangeably but there are differences between them.

- Opinion: A conclusion open to dispute (because different experts have different opinions)
- View: subjective opinion
- Belief: deliberate acceptance and intellectual assent
- Sentiment: opinion representing one’s feelings

II. LITERATURE SURVEY

The bag-of- words model is one of the most widely used feature model for almost all text classification tasks due to its simplicity coupled with good performance. The model represents the text to be classified as a bag or collection of individual words with no link or dependence of one word with the other, i.e. it completely disregards grammar and order of words within the text. This model is also very popular in sentiment analysis and has been used by various researchers. The simplest way to incorporate this model in our classifier is by using unigrams as features. Generally speaking n-grams is a contiguous sequence of “n” words in our text, which is completely independent of any other words or grams in the text. So unigrams is just a collection of individual words in the text to be classified, and we assume that the probability of occurrence of one word will not be affected by the presence or absence of any other word in the text. This is a very simplifying assumption but it has been shown to provide rather good performance.. One simple way to use unigrams as features is to assign them with a certain prior polarity, and take the average of the overall polarity of the text, where the overall polarity of the text could simply be calculated by summing the prior polarities of

individual unigrams. Prior polarity of the word would be positive if the word is generally used as an indication of positivity, for example the word “sweet”; while it would be negative if the word is generally associated with negative connotations, for example “evil”. There can also be degrees of polarity in the model, which means how much indicative is that word for that particular class. A word like “awesome” would probably have strong subjective polarity along with positivity, while the word “decent” would although have positive prior polarity but probably with weak subjectivity. There are three ways of using prior polarity of words as features. The simpler un-supervised approach is to use publicly available online lexicons/dictionaries which map a word to its prior polarity. The Multi-Perspective- Question-Answering (MPQA) is an online resource with such a subjectivity lexicon which maps a total of 4,850 words according to whether they are “positive” or “negative” and whether they have “strong” or “weak” subjectivity. The SentiWordNet 3.0 is another such resource which gives probability of each word belonging to positive, negative and neutral classes. The second approach is to construct a custom prior polarity dictionary from our training data according to the occurrence of each word in each particular class. For example if a certain word is occurring more often in the positive labeled phrases in our training dataset (as compared to other classes) then we can calculate the probability of that word belonging to positive class to be higher than the probability of occurring in any other class

$$(\textit{phrase}) = \log_2 \frac{\textit{hits}(\textit{phrase NEAR} \textit{“excellent”})}{\textit{hits}(\textit{“excellent”})}$$

$$(\textit{phrase NEAR} \textit{“poor”}) / (\textit{“excellent”})$$

Where $\textit{hits}(\textit{phrase NEAR} \textit{“excellent”})$ means the number of documents returned by the search engine in which the phrase (whose polarity is to be calculated) and word “excellent” are co-occurring. While $\textit{hits}(\textit{“excellent”})$ means the number of documents returned which contain the word “excellent”. Prabowo et al. have gone ahead with this idea and used a seed of 120 positive words and 120 negative to perform the internet searches. So the overall semantic orientation of the word under consideration can be found by calculating the closeness of that word with each one of the seed words and taking an average of it. Another graphical way of calculating polarity of adjectives has been discussed by **Hatzivassiloglou et al.** The process involves first identifying all conjunctions of adjectives from the corpus and using a supervised algorithm to mark every pair of adjectives as belonging to the same semantic orientation or different. A graph is constructed in which the nodes are the adjectives and links indicate same or different semantic orientation. Finally a clustering algorithm is applied which divides the graph into two subsets such that nodes within a subset mainly contain links of same orientation and links between the two subsets mainly contain links of different orientation. One of the subsets would contain positive adjectives and the other would contain negative.

Lina Zhou et al., investigated movie review mining using machine learning and semantic orientation. Supervised classification and text classification techniques are used in the proposed machine learning approach to classify the movie review. A corpus is formed to represent the data in the documents and all the classifiers are trained using this corpus. Thus, the proposed technique is more efficient. Though, the machine learning approach uses supervised learning, the proposed semantic orientation approach uses “unsupervised learning” because it does not require prior training in order to mine the data. Experimental results showed that the supervised approach achieved 84.49% accuracy in three-fold cross validation and 66.27% accuracy on hold-out samples. The proposed semantic orientation approach achieved 77% accuracy of movie reviews. Thus, the study concludes that the supervised machine learning is more efficient but requires a considerable amount of time to train the model. On the other hand, the semantic orientation approach is slightly less accurate but is more efficient to use in real time applications. The results confirm that it is practicable to automatically mine opinions from unstructured data.

Bo Pang et al., used machine learning techniques to investigate the effectiveness of classification of documents by overall sentiment. Experiments demonstrated that the machine learning techniques are better than human produced baseline for sentiment analysis on movie review data. The experimental setup consists of movie-review corpus with randomly selected 700 positive sentiment and 700 negative sentiment reviews. Features based on unigrams and bigrams are used for classification. Learning methods Naïve Bayes, maximum entropy classification and support vector machines were employed. Inferences made by Pang et al., is that machine learning techniques are better than human baselines for sentiment classification. Whereas the accuracy achieved in sentiment classification is much lower when compared to topic based categorization.

Zhu et al., proposed aspect based opinion polling from free form textual customer reviews. The aspect related terms used for aspect identification was learnt using a multi-aspect bootstrapping method. A proposed aspect-based segmentation model, segments the multi aspect sentence into single aspect units which was used for opinion polling. Using an opinion polling algorithm, they tested on real Chinese restaurant reviews achieving 75.5 percent accuracy in aspect-based opinion polling tasks. This method is easy to implement and is applicable to other domains like product or movie reviews.

Jeonghee Yi et al., proposed a Sentiment Analyzer to extract opinions about a subject from online data documents. Sentiment analyzer uses natural language processing techniques. The Sentiment analyzer finds out all the references on the subject and sentiment polarity of each reference is determined. The sentiment analysis conducted by the researchers utilized the sentiment lexicon and sentiment pattern database for extraction and association purposes. Online product review articles for digital camera and music were analyzed using the system with good results.

I. PROBLEM DEFINITION

Searching problem

- We have to find a particular word in about 2500 files.
- All words are weighted same for example good and best belongs to same Category.

- The sequence in which words come in test data is neglected.
- Other issues – Efficiency provided from this implementation is only 40-50%.

Existing System

In the existing system with the evolution of web technology, there is a huge amount of data present in the web for the internet users. These users not only use the available resources in the web, but also give their reviews, thus generating additional useful information. Due to overwhelming amount of user's reviews available through the web resources but using these reviews we cannot find which changes are required for a better college system. In this existing system we will get only reviews from users but not analyze the reviews.

Disadvantages of existing system

- Not get all reviews from other websites.
- There are not reviews analysis only get reviews from user.
- Cannot idea for which changes are required for system.

II. PROPOSED SOLUTION

To overcome the existing drawbacks we can propose a system that can extract the college reviews information from other websites using web mining technique and analyze these reviews using sentiment analysis. Sentiment analysis is a very relevant technique nowadays for analysis. Sentiment analysis or web mining is the process of automatically extracting knowledge from sentiments or reviews of others about some topic or problem. We can identify reviews in a large unstructured/structured data and analyze the polarity of reviews.

In this proposed system we can use Naive Bayes algorithm for analyzing college reviews and to tag a given review as positive or negative. The results can be used for various purposes such as guiding decisions to improve the college system.

Advantages of Proposed system

1. Easily gets reviews from various college websites.
2. Sentiment analysis gives a proper result of positive or negative reviews.
3. Using this analysis we can easily get what is our system's plus point and which sectors require changes.
4. Gives better services for user.

III. METHODOLOGY

Naive-Bayes Classification Algorithm

The Bayesian Classification represents a supervised learning method as well as a statistical method for classification. It assumes an underlying probabilistic model and it allows us to capture uncertainty about the model in a principled way by determining probabilities of the outcomes. It can solve diagnostic and predictive problems.

Bayesian classification provides practical learning algorithms and prior knowledge and observed data can be combined. Bayesian Classification provides a useful perspective for understanding and evaluating many learning algorithms. It calculates explicit probabilities for hypothesis and it is robust to noise in input data.

1. Dictionary Generation:- Count occurrence of all words in our whole data set and make a dictionary of some most frequent words.
2. Feature set Generation- All documents are represented as a feature vector over the space of dictionary words.
- For each document, keep track of dictionary words along with their number of occurrence in that document. Calculate Probability of occurrence of each label. Here label is negative and positive.
3. Training :- In this phase we have to generate training data (words with probability of occurrence in positive/negative training datafiles). Calculate for each label. Calculate for each dictionary word and store the result (Here: label will be negative and positive). Now we have word and corresponding probability for each of the defined labels.

Module Description

1. Login - In this module using username and password user logs into the system. In this login system authentication of user so only valid person logs into the system.

2. Data Collection - In this user select one college name from a college list and click on submit after submitting our system get reviews of this college using web mining technique. In web mining technique the system get data from another websites where college reviews are present related with this college.
3. Sentiment Analysis - The reviews retrieved using web mining technique from another websites which can be analyzed using Naive Bayes algorithm and get result as positive and negative reviews.

Output

In this module we can display out for user. The output display college information, Placement, Teaching or faculty, Crowd and display pie chart of positive and negative percentage and maximum 10 comments.

IV. REQUIREMENTS

Hardware:

1. Processor: Pentium 4 , RAM: 2GB or more, Hard disk: 16 GB or more

Software:

1. Apache tomcat Server, Windows Operating System, Eclipse, Java, MYSQL

V. FUTURE ASPECTS

Here this kind of a system is in the developing phase and lots of future enhancements are planned and are undergoing 1st level analysis. This application can be expanded with many new other building schemes and areas. Due to the time constraint we were not able to provide various enhancements such as:

This is a Reference Model of share trading Application for Education Institute. You can add following Futures. Change the share price directly through internet.

- Login only those users there have account in any bank system.
- Admin can put various advertisements

REFERENCES

- [1] Bing Liu, "Exploring User Opinions in Recommender Systems", Proceeding of the second KDD workshop on Large Scale Recommender Systems and the Netflix Prize Competition, Aug 24, 2008, Las Vegas, Nevada, USA.
- [2] Dave D., Lawrence A., Pennock D., "Mining the Peanut Gallery: Opinion Extraction and Semantic Classification of Product Reviews", Proceedings of International World Wide Web Conference (WWW'03), 2003.
- [3] Turney, P., "Thumbs Up or Thumbs Down? Semantic Orientation Applied to Unsupervised Classification of Reviews", ACL'02, 2002.
- [4] Lina Zhou, Pimwadee Chaovalit, "Movie Review Mining: a Comparison between Supervised and Unsupervised Classification Approaches", Proceedings of the 38th Hawaii International Conference on System Sciences, 2005.
- [5] Bo Pang, Lillian Lee, and Shivakumar Vaithyanathan, "Thumbs up? Sentiment classification using machine learning techniques", In Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 79-86, 2002.
- [6] Zhu, Jingbo Wang, Huizhen Zhu, Muhua Tsou, Benjamin K. Ma, Matthew, "Aspect-Based Opinion Polling from Customer Reviews", IEEE Transactions on Affective Computing, Volume: 2, Issue: 1 On page(s): 37. Jan-June 2011.40
- [7] Yi, J., T. Nasukawa, R. Bunescu, and W. Niblack: 2003, "Sentiment Analyzer: Extracting Sentiments about a Given Topic using Natural Language Processing Techniques", In: Proceedings of the 3rd IEEE International Conference on Data Mining (ICDM-2003). Melbourne Florida.