

## A Review on Subjectivity Analysis through Text Classification Using Mining Techniques

Ashwin Shinde\*, Mayuri Marawar\*\*, Minal Domke\*\*\*, Bhushan Manjre\*\*\*\*

\*(Department of Computer Science and Engineering, RTMNU University, Nagpur, Maharashtra)

\*\*(Department of Computer Science and Engineering, RTMNU University, Nagpur, Maharashtra)

\*\*\* (Department of Information Technology, RTMNU University, Nagpur, Maharashtra)

\*\*\*\*(Department of Information Technology, RTMNU University, Nagpur, Maharashtra)

### ABSTRACT

The increased use of web for expressing ones opinion has resulted in to an enhanced amount of subjective content available in the Web. These contents can often be categorized as social content like movie or product reviews, Customer Feedbacks, Blogs, Communication exchange in discussion forums etc. Accurate recognition of the subjective or sentimental web content has a number of benefits. Understanding of the sentiments of human masses towards different entities and products enables better services for contextual advertisements, recommendation systems and analysis of market trends. The objective behind framing this paper to analyze various sentiment based classification techniques which can be utilized for quick estimation of subjective contents of Political reviews based on politicians speech. The paper elaborately discusses supervised machine learning algorithm: Naïve Bayes classification and compares its overall accuracy, precisions as well as recall values.

**Keywords:** Sentiment Analysis, Text Classification, Probabilistic Estimation.

### I. INTRODUCTION

Data mining is a process of mined valuable data from a large set of data. Several analysis tools of data mining (like, clustering, classification, regression etc,) can be used for sentiment analysis task [1][2].

Sentiment analysis is critical because helps you see what customers like and dislike about you and your brand. Customer feedback—from social media, your website, your call center agents, or any other source—contains a treasure trove of useful business information. But, it isn't enough to know what customers are talking about. You must also know how they feel. Sentiment analysis is one way to uncover those feelings.

Sometimes known as “opinion mining,” sentiment analysis can let you know if there has been a change in public opinion toward any aspect of your business. Peaks or valleys in sentiment scores give you a place to start if you want to make product improvements, train sales or customer care agents, or create new marketing campaigns.

Sentiment analysis is not a once and done effort. By reviewing your customer's feedback on your business regularly you can be more proactive regarding the changing dynamics in the market place.

Sentiment mining is one of the important aspects of data mining. Subjective analysis is the detection of attitudes “enduring, affectively colored beliefs, dispositions towards objects or persons”. In

Sentiment mining important data can be mined based on the positive or negative senses of the collected data. Sentiment Analysis also known as Opinion Mining refers to the use of natural language processing, text analysis and computational linguistic to identify and extract subjective information in source materials. Here the source materials refer to opinions / reviews /comments given in various social networking sites [3].The Sentiment found within comments, feedback or critiques provide useful indicators for many different purposes and can be categorized by polarity [4].By polarity we tend to find out if a review is overall a positive one or a negative one. For example: 1) Positive Sentiment in subjective sentence: “I loved the movie ”—This sentence is expressed positive sentiment about the movie and we can decide that from the sentiment threshold value of word “loved”. So, threshold value of word “loved” has positive numerical threshold value. 2) Negative sentiment in subjective sentences: “as worst as it can get” or “Flop flick” defined sentence is expressed negative sentiment about the movie named “Flop flick” and we can decide that from the sentiment threshold value of word “Flop”. So, threshold value of word “Flop” has negative numerical threshold value. Sentiment Analysis is of three different types: Document level, Sentence level and Entity level. The traditional text mining concentrates on analysis of facts whereas opinion mining deals with the attitudes [5].

NLP, statistics, or machine learning methods are used to extract, identify, or otherwise characterize the sentiment content of a text unit. Sometimes referred to as opinion mining, although the emphasis in this case is on extraction.

## II. NAÏVE BAYESIAN CLASSIFICATION

“What are Bayesian classifiers?” Bayesian classifiers are statistical classifiers. They can predict class membership probabilities such as the probability that a given tuple belongs to a particular class. Bayesian classification is based on Bayes’ theorem, described next. Studies comparing classification algorithms have found a simple Bayesian classifier known as the naïve Bayesian classifier to be comparable in performance with decision tree and selected neural network classifiers. Bayesian classifiers have also exhibited high accuracy and speed when applied to large databases. Naïve Bayesian classifiers assume that the effect of an attribute value on a given class is independent of the values of the other attributes. This assumption is called class conditional independence. It is made to simplify the computations involved and, in this sense, is considered “naïve.”

The **naïve Bayesian** classifier, or **simple Bayesian** classifier, works as follows:

1. Let  $D$  be a training set of tuples and their associated class labels. As usual, each tuple is represented by an  $n$ -dimensional attribute vector,  $\mathbf{X} = [x_1, x_2, \dots, x_n]$ , depicting  $n$  measurements made on the tuple from  $n$  attributes, respectively,  $A_1, A_2, \dots, A_n$ .

2. Suppose that there are  $m$  classes,  $C_1, C_2, \dots, C_m$ . Given a tuple,  $\mathbf{X}$ , the classifier will predict that  $\mathbf{X}$  belongs to the class having the highest posterior probability, conditioned on  $\mathbf{X}$ . That is, the naïve Bayesian classifier predicts that tuple  $\mathbf{X}$  belongs to the class  $C_i$  if and only if

$$P(C_i|\mathbf{X}) > P(C_j|\mathbf{X})$$

Thus, we maximize  $P(C_i|\mathbf{X})$ . The class  $C_i$  for which  $P(C_i|\mathbf{X})$  is maximized is called the maximum posterior hypothesis. By Bayes’ theorem

$$P(C_i|\mathbf{X}) = P(\mathbf{X}|C_i)P(C_i)/P(\mathbf{X})$$

3. As  $P(\mathbf{X})$  is constant for all classes, only  $P(\mathbf{X}|C_i)P(C_i)$  needs to be maximized. If the class prior probabilities are not known, then it is commonly assumed that the classes are equally likely, that is,  $P(C_1) = P(C_2) = \dots = P(C_m)$ , and we would therefore maximize  $P(\mathbf{X}|C_i)$ . Otherwise, we maximize  $P(\mathbf{X}|C_i)P(C_i)$ . Note that the class prior probabilities may be estimated by  $P(C_i) = |C_i|/|D|$ , where  $|C_i|$  is the number of training tuples of class  $C_i$  in  $D$ .

4. Given data sets with many attributes, it would be extremely computationally expensive to compute  $P(\mathbf{X}|C_i)$ . To reduce computation in evaluating  $P(\mathbf{X}|C_i)$ , the naïve assumption of class-conditional

independence is made. This presumes that the attributes’ values are conditionally independent of one another, given the class label of the tuple (i.e., that there are no dependence relationships among the attributes). Thus,

$$P(\mathbf{X}|C_i) = P(x_1|C_i)x_1 P(x_2|C_i)x_2 \dots x_n P(x_n|C_i)$$

We can easily estimate the probabilities  $P(x_1|C_i), P(x_2|C_i), \dots, P(x_n|C_i)$  from the training tuples. Recall that here  $x_k$  refers to the value of attribute  $A_k$  for tuple  $X$ . For each attribute, we look at whether the attribute is categorical or continuous-valued. Now if we want to make sentiment analysis of people’s reviews regarding the speech made available on some blog of a politician, then we can classify the sentiment by using the naïve bayes classifier as below.

First we need the training data which is represented in tabular form as below.

**Table 1.** Training Data

Name	Age	Work done	Genuineness in speech	Potential Leader
Amit	22	A bit	High	Yes
Ajay	30	Moderate	High	Yes
Mohit	55	Moderate	Low	No
Raj	70	High	Low	No
Ravi	45	Moderate	High	Yes

This training data will then be fed to the classification algorithm and then as per the classification rules inferred from the training data classification of good leader or bad leader will be done. Classification rules can be given as below.

IF age\_group=20-35 , work done = a Bit and genuineness=High THEN descigion column or class label column will yield result as yes for potential leadership . IF age\_group>60, work done=High and genuineness=Low THEN descigion column or class label column will yield No for potential Leadership. Now suppose if the new entry in the table is like age=72 and genuineness=Low THEN decision column or class label column will yield result as Potential Leader.

We wish to predict the class label of a tuple using naïve Bayesian classification, given the same training data as given in the table . The data tuples are described by the attributes Name, age, work done, genuineness. The class label attribute, Potential Leader, has two distinct values (namely, yes, no). Let  $C_1$  correspond to the class Potential Leader as yes and  $C_2$  correspond to Potential Leader as no. The tuple we wish to classify is

$$\mathbf{X} = (\text{age} = 23, \text{work done}=a \text{ bit}, \text{genuineness}=high)$$

We need to maximize  $P(\mathbf{X}|C_i)P(C_i)$ , for  $i = 1, 2$ .  $P(C_i)$ , the prior probability of each class, can be computed based on the training tuples:

$P(\text{Potential Leader} = \text{yes}) = 3/5 = 0.6$   
 $P(\text{Potential Leader} = \text{No}) = 2/5 = 0.4$   
To compute  $P(\mathbf{X}|C_i)P(C_i)$ , for  $i = 1, 2$ .  $P(C_i)$  we compute following conditional probabilities.  
 $P(\text{Age}=20-35|\text{Potential Leader}=\text{Yes})=2/3=0.6$   
 $P(\text{Age}=20-35|\text{Potential Leader}=\text{No})=0/3=0$   
 $P(\text{Work Done}=\text{a bit}|\text{Potential Leader}=\text{Yes})=1/3=0.3$   
 $P(\text{Work Done}=\text{a bit}|\text{Potential Leader}=\text{No})=0/3=0$   
 $P(\text{Genuineness}=\text{High}|\text{Potential Leader}=\text{Yes})=3/3=1$   
 $P(\text{Genuineness}=\text{Low}|\text{Potential Leader}=\text{No})=2/3=0.6$ .  
Using these probabilities we obtain:  
 $P(\mathbf{X}|\text{Potential Leader}=\text{yes})=P(\text{Age}=20-35|\text{Potential Leader}=\text{Yes}) \times P(\text{Work Done}=\text{a bit}|\text{Potential Leader}=\text{Yes}) \times P(\text{Genuineness}=\text{High}|\text{Potential Leader}=\text{Yes})=0.6 \times 0.3 \times 1=0.18$ . and  $P(\mathbf{X}|\text{Potential Leader}=\text{No})=P(\text{Age}=20-35|\text{Potential Leader}=\text{No}) \times P(\text{Work Done}=\text{a bit}|\text{Potential Leader}=\text{No}) \times P(\text{Genuineness}=\text{Low}|\text{Potential Leader}=\text{No})=0 \times 0 \times 0.6=0$   
To find the class,  $C_i$ , that maximizes  $P(\mathbf{X}|C_i)P(C_i)$ , we compute:  
 $P(\mathbf{X}|\text{Potential Leader}=\text{yes}) \quad P(\text{Potential Leader}=\text{yes})=0.18 \times 0.6=0.1$   
 $P(\mathbf{X}|\text{Potential Leader}=\text{no}) \quad P(\text{Potential Leader}=\text{no})=0 \times 0.4=0$   
Therefore, the naive Bayesian classifier predicts Potential Leader= yes for tuple  $\mathbf{X}$ .

### III. CONCLUSION

The aim of study is to evaluate the performance for sentiment classification in terms of accuracy, precision and recall. In this paper, we presented supervised machine learning algorithms of Naive Bayes for sentiment classification of the reviews of peoples from reading political leaders speech recorded in textual form on blog. The approach show that the classifiers yields better result for the with the Naive Bayes approach giving accuracies and. Thus we can say Naïve Bayes classifier can be used successfully to analyse political reviews.

### REFERENCES

- [1] L. Dey and S. Chakraborty, "Canonical PSO Based K-Means Clustering Approach for Real Datasets", International Scholarly Research Notices, Hindawi Publishing Corporation, Vol.2014,pp.1- 11,2014.
- [2] R. Dey and S. Chakraborty, "Convex-hull & DBSCAN clustering to predict future weather", 6th International IEEE Conference and Workshop on Computing and Communication, Canada, 2015, pp.1- 8.
- [3] Lina L. Dhande and Dr. Prof. Girish K. Patnaik, "Analyzing Sentiment of Movie Review Data using Naive Bayes Neural Classifier", IJETTCS, Volume 3, Issue 4 July-August 2014, ISSN 2278-6856.
- [4] P.Kalaivani, "Sentiment Classification of Movie Reviews by supervised machine learning approaches" et.al, Indian Journal of Computer Science and Engineering (IJCSE) ISSN : 0976-5166 Vol. 4 No.4 Aug-Sep 2013.
- [5] Meena Rambocas, João Gama, "Marketing Research: The Role of Sentiment Analysis", April 2013, ISSN: 0870-8541
- Chapters in Books:**
- [6] Jiawei Han, Michelin Kamber, Data Mining Concepts and Techniques, Mourgan Kaupmann ,Third Edition,351-354.