# Real Time Sentiment Classification Using Unsupervised Reviews

E.Divya

**Abstract-** Sentiment classification is an important task in everyday life. Users express their opinion about their product, movies and so on. All the web page contains reviews that are given by users expressing different polarity i.e. positive or negative. It is useful for both the producer and consumer to know what people think about the particular product or services based on their reviews. Automatic document classification [2],[3] is the task of classifying the reviews based on the sentiment expressed by the reviews. Sentiment is expressed differently in different domains. The data trained on one domain cannot be applied to the data trained on another domain [6]. The cross domain sentiment classification overcomes these problems by creating thesaurus for labeled data on the target domain and unlabeled data from source and target domains. Sentiment sensitivity is achieved by creating thesaurus. The sentiments cannot be analyzed for sentence and the data to be trained on a particular domain. The proposed method focus on unsupervised method of using Twitter, the most popular micro blogging platform, for the task of Opinion analysis. A sentiment classifier, that is able to determine positive, negative and neutral sentiments for a twitter website reviews.

**Key terms -** Sentiment Analysis, Thesaurus, Sentiment Classification.

————————— ◆ —————————

## 1 INTRODUCTION

### 1.1 SENTIMENT ANALYSIS

Sentiment analysis is used in natural language processing. Its main aim is to identify and extract sensitive information in the source. Sentiment analysis is a recent attempt to deal with evaluative aspects of text. In sentiment analysis, one fundamental problem is to recognize whether given text expresses positive or negative evaluation. Such property of text is called polarity. Sentiment classification can be applied in various tasks such as opinion mining [4], opinion summarization [5], contextual advertising [6] and market analysis [7].

Supervised learning technique that requires labeled data have been successfully used for building sentiment classifier for particular domain [2]. Supervised learning is the machine learning task of inferring a function from labeled trained data. The training data set consist of a set of training examples. In supervised learning, each data set is a pair consisting of an input object (typically a vector) and a desired output value. A supervised learning algorithm analyzes the training data and produces a conditional

_____

• *E.Divya is currently pursuing masters degree program in Computer science and Engineering in Sri Krishna College of Engineering and Technology.Email_id : divyaeswarancse@gmail.com*

function, which can be used for mapping new data. An optimal scheme will allow for the algorithm to correctly determine the class labels for unseen instances.

Unsupervised Learning technique is used for classify the review as recommended or not recommended [2].It is used to find the hidden data from the unlabeled data. The algorithm takes the review as input and gives a classification as output. The features and working of sentiment classifications are done by various level of sentiment analysis.

### 1.2 DOCUMENT LEVEL SENTIMENT ANALYSIS

This is the simplest form of sentiment analysis and it is assumed that the document contains an opinion on one main object expressed by the author of the document. There are two main approaches to document-level sentiment analysis: supervised learning and unsupervised learning. The supervised approach assumes that there is a finite set of classes into which the document should be classified and training data is available for each class that is positive and negative. Simple extensions can also added a neutral class. With the training data, the system learns a classification model by using one of the common classification algorithms such as SVM [2], Naïve Bayes[2],Turney[3].This

classification is then used to tag new documents into their various sentiment classes. When a numeric value (in some finite range) is to be assigned to the document then regression can be used to predict the value to be assigned to the document.

## 1.3 SENTENCE-LEVEL SENTIMENT ANALYSIS

A single document may contain multiple opinions even about the same data. When we want to have a more detailed view of the different opinions expressed in the document about the entities we must move to the sentence level. Before analyzing the polarity of the sentences we must determine if the sentences are subjective or objective. Only subjective sentences will be further analyzed. After we have zoned in on the subjective sentences we can classify these sentences into positive or negative classes. Sentence-level sentiment analyses are either based on supervised learning or on unsupervised learning.

## 1.4 ASPECT BASED SENTIMENT ANALYSIS

Aspect-based sentiment analysis (also called feature-based sentiment analysis) is the research problem that focuses on the recognition of all sentiment expressions within a given document and the aspects to which they refer.

The classic approach, which is used by many commercial applications, to the identification of all aspects in a corpus of product reviews is to extract all noun phrases (NPs) and then keep just the NPs whose frequency is above some experimentally determined threshold. One approach is to reduce the noise in the found NPs. The main idea is to measure for each candidate NP the PMI[3] with phrases that are tightly related to the product category

## 2 RELATED WORKS

Sentiment analysis (also known as opinion mining) refers to the use of natural language processing, text analysis and computational linguistics to identify and extract subjective information in source materials. Sentiment is expressed differently in different domains. The data trained on one domain cannot be applied to the data trained on another domain. The cross domain sentiment classification overcomes these problems by creating thesaurus [1] for labeled data on the target domain and unlabeled data from source and target domains. Sentiment sensitivity is achieved by creating thesaurus. The created thesaurus is used to expand the feature vector. Reviews are taken from different products and the thesaurus is created for multiple domains which contain both positive and negative words. Thus the created sentiment sensitive thesaurus[1] captures the words with similar sentiment.

A cross domain sentiment classification system must overcome two main challenges [8],[9]. First, we must identify which source domain features are related to which target domain features. Second, we require a learning framework to incorporate the information regarding the relatedness of source and target domain features.

The problem as one of feature expansion, where we append additional related features to feature vectors that represent source and target domain reviews to reduce the mismatch of features between the two domains. Methods that use related features have been successfully used in numerous tasks such as query expansion [10], information retrieval [11], and document classification [12].

We create a sentiment sensitive thesaurus that aligns different words that express the same sentiment in different domains. We use labeled data from multiple source domains and unlabeled data from source and target domains to represent the distribution of features. We use lexical elements (unigrams and bigrams of word lemma) and sentiment elements (rating information) to represent a user review. Next, for each lexical element we measure its relatedness to other lexical elements and group related lexical elements to create a sentiment sensitive thesaurus. The thesaurus captures the relatedness among lexical elements that appear in source and target domains based on the contexts in which the lexical elements appear (its distributional context).

## 2.1 Thesaurus

A thesaurus is a reference work that lists words grouped together according to similarity of meaning (containing synonyms and sometimes antonyms), in contrast to a dictionary, which provides definitions for words, and generally lists them in alphabetical order. The main purpose of such reference works is to help the user to find the words.

The thesaurus is constructed by calculating the point wise mutual information between a lexical element w and the feature vector v. Point wise mutual information between a lexical element and co occurrences word.

$$f(u, w) = \log\left(\frac{\frac{c(u,w)}{N}}{\frac{\sum_{i-1}^{n} c(i,w)}{N} \times \frac{\sum_{j-1}^{m} c(u,j)}{N}}\right) \qquad (3.1)$$

**C(u,w)** denotes the number of review sentences in which a lexical element u and a feature w co-occur, n and m denote the total number of lexical elements and the total number of features.
To compute the relatedness of the element

$$\tau(v, u) = \frac{\sum_{\omega \in \{x|f(v,x)>0\}} f(u,w)}{\sum_{\omega \in \{x|f(v,x)>0\}} f(u,w)} \qquad (3.2)$$

The relatedness score $\tau(v,u)$ can be interpreted as the proportion of pmi-weighted features of the element u that
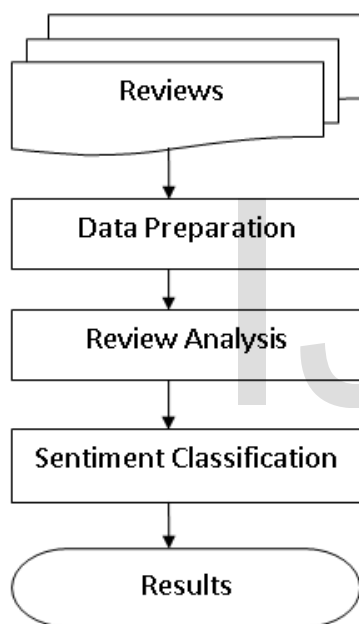
are shared with element v.

# 3 PROPOSED WORK

## 3.1 Opinion Mining Process

Opinion mining (or sentiment analysis) has attracted great interest in recent years, both in academia and industry due to its potential applications. One of the most promising applications is analysis of opinions in social networks. Lots of people write their opinions in forums, micro blogging or review websites. This data is very useful for business companies, governments, and individuals, who want to track automatically attitudes and feelings in those sites. Namely, there is a lot of data available that contains much useful information, so it can be analyzed automatically.

FIG .1 SENTIMENT ANALYSIS MODEL



For instance, a customer who wants to buy a product usually searches the Web trying to find opinions of other customers or reviewers about this product. In fact, these kinds of reviews affect customer's decision. Twitter is used for the experiments. Twitter is a micro blogging platform where users post their messages, opinions, comments, etc. Contents of the messages range from personal thoughts to public statements. A Twitter message is known as tweet. Tweets are very short; the maximum size of a tweet is 140 characters that usually correspond to a phrase. Thus, the work is limited to sentence level.

## 3.1.1 PREPROCESSING

Analysis of tweets is complex task because these messages

are full of slang, misspellings [7] and words borrowed from other languages. So in order to tackle the problems mentioned and to deal with the noise in texts, we normalize the tweets before training the classifiers.

## 3.2 Definitions
### POS TAGGER

Part-of-speech tagging (POS tagging or POST),[13] also called grammatical tagging or word-category disambiguation, is the process of marking up a word in a text. It is also called Tokenization.

### LEMMATIZATION

Lemmatization is the process of grouping together the infected form of feature into a single item. Lemmatization is similar to stemming.

### STOP WORD

Stop word is the process of filtering out the prior to or after preprocessing.

## 3.2 SELECTED CLASSIFIERS

Existing method uses various machine learning classifiers. The machine learning classifiers we selected were: Naïve Bayes (NB), Turney classifier

## 3.3 CLASSIFIERS
### 3.3.1 Bayesian Opinion Mining

The Naive Bayesian classier is a straightforward and frequently used method for supervised learning [2]. It provides a flexible way for dealing with any number of attributes or classes, and is based on probability theory. It is the asymptotically fastest learning algorithm that examines all its training input.

Bayesian classifiers are based around the Bayes rule, a way of looking at conditional probabilities that allows you to flip the condition around in a convenient way. A conditional probably is a probably that event X will occur, given the evidence Y. That is normally written $P(X \mid Y)$. The Bayes rule allows us to determine this probability when all we have is the probability of the opposite result, and of the two components individually: $P(X \mid Y) = P(X)P(Y \mid X) / P(Y)$. This restatement can be very helpful when we're trying to estimate the probability of something based on examples of it occurring.

Formula looks like this.

**P(sentiment|sentence)=P(sentiment)P(sentence sentiment)/P(sentence)**

In this case, we're trying to estimate the probability that a document is positive or negative, given its contents. We can restate that so that is in terms of the probability of that document occurring if it has been predetermined to be

positive or negative. The thing that makes this a "naive" Bayesian process is that we make a big assumption about how we can calculate at the probability of the document occurring: that it is equal to the product of the probabilities of each word within it occurring. This implies that there is no link between one word and another word. This independence assumption is clearly not true: there are lots of words which occur together more frequently that either does individually, or with other words, but this convenient fiction massively simplifies things for us, and makes it straightforward to build a classifier.

We can estimate the probability of a word occurring given a positive or negative sentiment by looking through a series of examples of positive and negative sentiments and counting how often it occurs in each class. This is what makes this supervised learning - the requirement for pre-classified examples to train on.

We can drop the dividing P(line), as it's the same for both classes, and we just want to rank them rather than calculate a precise probability. We can use the independence assumption to let us treat P(sentence | sentiment) as the product of P( token | sentiment) across all the tokens in the sentence.

**So, we estimate P (token | sentiment) as Count (this token in class) + 1 / count (all tokens in class) + count (all tokens)**

The extra 1 and count of all tokens is called 'add one' or Laplace smoothing, and stops a 0 finding its way into the multiplications. If we didn't have it any sentence with an unseen token in it would score zero.

The classify function starts by calculating the prior probability (the chance of it being one or the other before any tokens are looked at) based on the number of positive and negative examples - in this example that'll always be 0.5 as we have the same amount of data for each. We then tokenize the incoming document, and for each class multiply together the likelihood of each word being seen in that class. We sort the final result, and return the highest scoring class.

### 3.2.2 Turney Opinion Mining
The Turney algorithm takes the input as reviews and produce classification as output. It contain three steps they are
**Algorithm**
The first step of the algorithm is to extract phrases containing adjectives or adverbs.
**Step 1:** Part-of-speech (POS) tagging Extracting two consecutive words (two word phrases) from reviews if their tags conform to some given patterns, e.g., (1) JJ, (2) NN.
The second step is to estimate the semantic orientation of the extracted phrases, using the PMI-IR algorithm [3]. This algorithm uses mutual information as a

measure of the strength of semantic association between two words
**Step 2:** Estimate the sentiment orientation (SO) of the extracted phrases

$$SO-A(word) = \sum_{pword \in Pwords} A(word, pword) - \sum_{nword \in Nwords} A(word, nword) \quad \textbf{(4.1)}$$

**Pwords** = a set of words with positive semantic orientation
**Nwords** = a set of words with negative semantic orientation
**A (word1, word2)** = a measure of association between word1 and word2
**Pwords**={good, nice, excellent, positive, fortunate, correct, and superior}
**Nwords**={bad, nasty, poor, negative, unfortunate, wrong, and inferior}.

**Positive Review:**
Example: love the local branch however communication may break down if they have to go through head office. Avg. SO Value=0.1414

**Negative Review:**
Example: Do not bank here, their website is even worse than their actual locations. Avg.So value: -0.0766

**Steps3:** Point wise Mutual Information (PMI), The Point wise Mutual Information (PMI) between two words, word1 and word2, is defined

$$PMI(word_1, word_2) \left( \frac{p(word_1 \& word_2)}{p(word_1)p(word_2)} \right) \textbf{(4.2)}$$

Here, p(word1 & word2) is the probability that word1 and word2 co-occur. If the words are statistically independent, the probability that they co-occur is given by the product p(word1) p(word2). The ratio between p(word1 & word2) and p(word1) p(word2) is a measure of the degree of statistical dependence between the words. The log of this ratio is the amount of [information that we acquire about the presence of one of the words when we observe the other.

## 5   CONCLUSION
The expression of opinions of users in specialized sites for evaluation of products and services, and also on social networking platforms, has become one of the main ways of communication, due to spectacular development of web environment in recent years. The large amount of information on these platforms make them viable for use as

data sources, in applications based on opinion mining and sentiment analysis. This paper presents a method of sentiment analysis, on the review made by users on twitter. Classification of reviews in both positive and negative classes is done based on a naive Bayes algorithm and Turney approach. This work focuses on twitter to collect real time opinions for the task of sentiment analysis. Evaluation is done to calculate the performance of the classifier.

# REFERENCES

[1] Danushka Bollegala,David Weir and John Carroll"Cross-Domain Sentiment Classification Using a Sentiment Sensitive Thesaurus",vol 25,No.8,pp.1719-1731,2013.

[2] B. Pang, L. Lee, and S. Vaithyanathan, "Thumbs Up? Sentiment Classification Using Machine Learning Techniques," Proc. ACL-02 Conf. Empirical Methods in Natural Language Processing (EMNLP '02), pp. 79-86, 2002.

[3] P.D. Turney, "Thumbs Up or Thumbs Down? Semantic Orientation Applied to Unsupervised Classification of Reviews," Proc. 40th Ann. Meeting on Assoc. for Computational Linguistics (ACL '02),pp. 417-424, 2002.

[4] B. Pang and L. Lee, "Opinion Mining and Sentiment Analysis,"Foundations and Trends in Information Retrieval, vol. 2, nos. 1/2,pp. 1-135, 2008.

[5] Y. Lu, C. Zhai, and N. Sundaresan, "Rated Aspect Summarization of Short Comments," Proc. 18th Int'l Conf. World Wide Web (WWW '09), pp. 131-140, 2009.

[6] T.-K. Fan and C.-H. Chang, "Sentiment-Oriented Contextual advertising, "Knowledge and Information Systems, vol. 23, no. 3,pp. 321-344, 2010.

[7] M.Hu and B.Liu,"Mining and Summarizing Customer Reviews",Proc.10t ACM SIGKDD Int'l Conf.Knowledge Discovery and Data Mining(KDD'04),pp.168-177,2004.

[8]J.Blitzer, M. Dredze, and F. Pereira, "Biographies, Bollywood,Boom-Boxes and Blenders: Domain Adaptation for Sentiment Classification," Proc. 45th Ann. Meeting of the Assoc. Computational Linguistics (ACL '07), pp. 440-447, 2007.

[9] S.J. Pan, X. Ni, J.-T. Sun, Q. Yang, and Z. Chen, "Cross-Domain Sentiment Classification via Spectral Feature Alignment," Proc.19th Int'l Conf. World Wide Web (WWW '10), 2010.

[10] H. Fang, "A Re-Examination of Query Expansion using Lexical Resources," Proc. Ann. Meeting of the Assoc. Computational Linguistics (ACL '08), pp. 139-147, 2008.

[11] G. Salton and C. Buckley, Introduction to Modern Information Retrieval. McGraw-Hill Book Company, 1983.

[12] D. Shen, J. Wu, B. Cao, J.-T. Sun, Q. Yang, Z. Chen, and Y. Li,"Exploiting Term Relationship to Boost Text Classification," Proc.18th ACM Conf. Information and Knowledge Management (CIKM '09),pp. 1637-1640, 2009.

[13] T. Briscoe, J. Carroll, and R. Watson, "The Second Release of the RASP System," Proc. COLING/ACL Interactive Presentation Sessions Conf., 2006.

[14] T. Joachim, "Text Categorization with Support Vector Machines: Learning with Many Relevant Features," Proc. 10th European Conf.Machine Learning (ECML '98), pp. 137-142, 1998.

[15] V. Hatzivassiloglou and K.R. McKeon, "Predicting the Semantic Orientation of Adjectives," Proc. Ann. Meeting of the Assoc. Computational Linguistics (ACL '97), pp. 174-181, 1997.

[16] J.M. Wiebe, "Learning Subjective Adjective from Corpora," Proc.17th Nat'l Conf. Artificial Intelligence and 12th Conf. Innovative Applications of Artificial Intelligence (AAAI '00), pp. 735-740, 2000.

1.

2.