# Aspect Based Sentiment Analysis using POS Tagging and TFIDF

**Kotagiri. Srividya, A.Mary Sowjanya**

*Abstract: Social media content on the internet is increasing day by day. Since media knowledge helps people in making decisions, web based businesses give their clients an opportunity to express their opinions about items available on the web in the form of surveys and reviews. Sentiment analysis can be used on product reviews or tweets, comments, blogs to infer individual's feelings or attitudes. Here Aspect Based Sentiment Analysis is used to extract most interesting aspect of a particular product from unlabeled text. We have developed two models for aspect/feature extraction.Model1 uses POS tagging whereas Model2 utilizes TFIDF .In Model 1 we start with noun phrase algorithm and extend it to adjectives and adverbs to extract all the aspect terms. In model2 after data preprocessing TDIDF technique is used. The relative importances of the aspects are calculated and the most important positive, negative and neutral aspects are presented to the user. Naïve Bayes, Support Vector machine, Decision Tree, KNN were used to classify the sentiment polarity of the generated aspects*.

*Index Terms*: *Aspects, Opinion Mining, Naïve Bayes, Support Vector Machine, KNN, Decision tree, Polarity*

## I. INTRODUCTION

As the web and its innovation, individuals gets an opportunity to communicate their perspectives, interests and assessments about the things they see or use in the form of audits and feedbacks. Organizations and item suppliers are keenly follow the client surveys and reviews on assess the opinions of their clients about their items. Since client's port their feedbacks over diverse websites service providers face challenges in finding general behavior or specific interest of their customers. Other category of mining is not enough for customer to make proper decisions. Also they cannot provide complete information of the product Accordingly looking top to bottom into aspects and entities has provided another guidance for research called aspect or feature based opinion investigation. Document level or sentence level isn't sufficient for customers to get an opinion on a product. Additionally they don't give total data of an item, for instance, a positive or negative survey of a specific item doesn't imply that the reviewer likes or dislikes all aspects of that item. Therefore looking in depth into aspects and entities has given way for  aspect or feature based sentiment analysis. If any individual whose requirements to purchase a mobile with magnificent camera quality can look just for the reviews about that aspect for example "picture quality" rather than overall reviews of that mobile.

Single aspect opinions are reviews where individuals concentrate just around one aspect of the item, while in multi-aspect, individuals express different opinions on more than one aspect in the same review. For example, "This book had a good storyline, but the paper quality is bad". Here, the reviewer gave a positive mention on the "story" aspect and negative mention on the "paper quality" aspect of the book. Dealing with multi aspect sentences is a difficult task. The important task in aspect based sentiment analysis are: Identifying the opinion words from a given sentence (Aspect extraction) and ordering those extracted aspects (Aspect Sentiment Classification).After getting the aspects based on their polarity, we apply classification algorithms like Navie Bayes, Support Vector Machine, Decision Tree, and KNN. Naive Bayes classification algorithms are the more scalable and   classify documents based on the frequency of words. They consider more number of features while training the classifier and follow Bayes theorem by assuming independence of its features. It is better classifier than Naive Bayes. A SVM Classifier is formally defined by a separating hyper plane, hence it can fairly separate the classes.

## II. BACKGROUND KNOWLEDGE

Zainuddin Et Al have proposed a hybrid approach for aspect based sentiment analysis which connects association rule mining, dependency parsing and Senti Word net to solve aspect-based sentiment analysis problem [1].Bing Liu's angle based opinion mining procedure has been to find customer inclinations about the travel industry items, especially hotels and restaurants. Tourism reviews surveys contain profitable data to utilize aspect based opinion mining approach [2]. A multi-aspect bootstrapping technique was proposed to learn aspect related terms of every aspect, and then an aspect based segmentation model was proposed to divide a multi perspective sentence into numerous single aspect unit [3]. D V Nagarjuna Devi et al,built a general procedure of 'Angle or Feature based Sentiment Analysis' by utilizing a classifier called Support Vector Machine (SVM).Initially aspects were separated from sentence level and opiniated words were removed, then Senti Word Net  and  SVM classifier were used[4]. An Incremental machine learning approach  was used  for examining thesentiments of the users by Rajalaxmi Et Al[5].Aspects were extracted from the given sentences dependening on the procedure of classifier and compared with the past strategy.The proposed Incremental Decision tree Classification (IDTC) strategy utilizedan iterative procedure to group the given inputs for the assessment of sentiment Classification. The results shown that SVM method is much better compared to bayes classification in all cases.

*Retrieval Number F7935088619/2019©BEIESP*
*DOI: 10.35940/ijeat.F7935.088619*
*Journal Website: www.ijeat.org*

1960

*Published By:*
*Blue Eyes Intelligence Engineering*
*& Sciences Publication*

A comparative analysis of sentiment orientation using SVM and Navie Bayes techniques was done on the movie data set[6]. To improve aspect extraction dimension aspect pruning was proposed[7].The Initial step is to calculate the frequency of each word and choosing the most frequent aspects then semantic likeness of non-frequent words is calculated and aspects which are not semantically identified with the product are eliminated. Two methods were proposed by Raisa and Jayasree [8] to find the general aspect categories address in review sentences.Kim Et Al[9] introduced an unsupervised strategy that applies association rule mining on co-occurrence frequency information from a corpus to find aspect categories. They have used a supervised variant that outperforms existing methods with an F1-score of 84%.Povoda Et Al have used SVM classifier to perform sentiment analysis in text documents, especially text valence detection. Their solution has been assessed with various languages – English, German, Czech and Spanish. The described learning is fully automatic, can be connected to any language and no complicated preprocessing is needed [10].

## III. PROPOSED SYSTEM

By using document and sentence level opinion mining, users are unable to get interestingness of any entity. Whereas in aspect based opinion mining by using polarity users get most interesting aspects. The goal is to find interesting aspects based upon polarity and relative importance. Positive, negative and neutral aspects will be find based on the relative importance. Finally Naive Bayes and Support Vector Machine are used to test the accuracy. We develop two models, one using POS tagging and the other using TF-IDF. These models classify opinions from reviews using aspect level classification to get positive, negative and neutral aspects based upon the polarity. The most interesting aspects will be extracted based upon their relative importance. Naive Bayes and Support Vector Machine are used to get the sentiment from the data.

**Model 1 (using POS tagging)**

### A. Stop word removal

Stop words are regular words that for the most part don't add to the significance of a sentence, in any event for the motivations behind data recovery and normal language handling. These are words, for example, 'the' and 'a'. Most web indexes will sift through stop words from hunt questions and records so as to spare space in their file. NLTK accompanies a stop words corpus that contains word records for some languages.

### B. Stemming

Stemming is a procedure to remove appends from a word, winding up with the stem. For instance, the stem of cooking is cook, and a good stemming calculation realizes that the "ing" addition can be removed. Stemming is most usually utilized via web crawlers for ordering words. Rather than putting away all types of a word, a web crawler can store just the stems, extraordinarily decreasing the measure of file while expanding recovery exactness. A standout amongst the most widely recognized stemming calculations is the Porter stemming calculation by Martin Porter. It is intended to remove and supplant understood additions of English words.

### C. Lemmatization

Lemmatization is fundamentally the same as stemming, yet is progressively much the same as equivalent word substitution. A lemma is a root word, rather than the root stem. So dissimilar to stemming, you are always left with a substantial word that implies something very similar. Be that as it may, the word you end up with can be completely different.

**Algorithm for Preprocessing**

Input: a text file containing all the reviews
Output: sentences free from stop words, and stemmed, lemmatized words.
Method: Tokenize the given data set D into sentences S1,S2,S3,…Si again tokenize the sentence into words W1,W2,W3,..Wj

For each word W in a sentence, S look for the presence of it in stop word. If you found it, skip that word else include it and go for another word.
Do: FOR EACH SENTENCE in Si DO

### D. Aspect based sentimental analysis

In aspect based sentimental analysis, the main terms to be extracted are the aspects. As it is a complicated process we consider only nouns as the aspect terms. Nouns can be extracted from the chunk grammar. Chunk grammar is different for different data sets.

### E. Parts of speech tagging

Grammatical feature labeling changes a sentence, a list of words, into a list of tuples, where each tuple is of the structure (word, tag). The tag is a grammatical feature which implies whether the word is a thing, modifier, action word, etc. Without the grammatical feature labels, a chunker can't realize how to extract phrases from a sentence. POS-Tagging is treated as the most salient part in sentiment analysis. It is vital to distinguish characteristics or highlights, conclusion words from review sentence. POS Tagging should be possible either static (physically) or with the assistance of POS-Tagger device. POS-Tagger labels every one of the words from text. We have utilize NLTK's POS tag was utilized to label each word in reviews. The grammar to extract noun phrases is as follows.

**Chunk Grammar**

```
(r"NP:{<NN.?>+<RB><JJ>|""<NN><RB><JJ>|""<JJ><NN
S><RB>|"
"<RB>+<JJ><NN><RB>|""<RB><NN.?>+<VB.?>*<NN.?
>+<RB>|"
"<IN><RB><JJ>|""<NN><JJ><RB>|""<RB><VB.?>+<JJ>
<NN>|"
 "<RB><NN.?>+<VB.?>+<RB>|"        "<JJ><NN><RB>*|"
"<RB>|"
"<JJ><NN><RB><JJ>}")
```

### F. Noun Phrases

Noun phrases are used in English language. we assumed that aspects can be found in noun phrases hence we identify noun phrases in the text for POS tagging in the form (Word ,Noun phrase) to accomplish this task. Algorithm for noun phases shows the algorithm to find noun phrases.

**Algorithm for POS-Tagging**

Input: POS-Tagged words
Output: a Tree with noun phrases
Method: Extract noun phrases for dataset D of Sentence Si according to the chunk grammar rule
Do: FOR EACH SENTENCE Si DO

### G. Build Word-net

WordNet is a lexical database for the English language. It is a word reference planned explicitly for natural language processing.

**Proposed algorithm for finding Aspect Polarity( Probabilities)**

Algorithm for Aspect Polarity
 Input: noun list, other list, word probabilities
Output: Aspect polarity as positive, negative and neutral
 Method: get noun list NL, other list OL
For each item in OL build wordnet ON
For each item in ON look in word probabilities
If True get probability and polarity
 Repeat
Sum up all probabilities is noun/ noun list
 If sum >=0.7 it is positive
 If sum <=0.4 it is negative
Else it is neutral Repeat

### H. Aspect Classification

 We used Naïve Bayes and SVM algorithms to verify the accuracy of aspect polarity generated above are applied on the extracted aspects to verify the accuracy of polarity classification done.

#### Model 2 (using TF-IDF)

Data preprocessing (Stop word removal, Stemming, Lemmatization) is done on the iPod data set. After which aspect based sentiment analysis is also carried out. Now we use TF-IDF feature extraction technique. TFIDF stands for Term Frequency Inverse Document Frequency and it is used in text mining and information retrieval. Term frequency measures the number of items a word appears in a document divided by the total number of words in that document. Whereas inverse document frequency is computed as the logarithm of the number of the documents in the corpus divided by the number of documents where the specific term appears. So TF gives frequency of a term and IDF gives importance of a term.

#### Algorithm

Input: sentences free from stop words, and stemmed, lemmatized words.
Output: Aspects that retrieved using tf-idf.
Method: Load the input data file.
Calculating the frequency of word in the given input file i.e., number of occurrences of i in j.

$TF(w) =$ Number of times term w appears in a document) / (Total number of terms in the document)
Number of documents containing i.

$IDF(w) = \log_e($Total number of documents / Number of documents with term w in it$)$

**TF-IDF=TF(w)*IDF(w)**

Now the algorithm proposed above for finding aspect polarities is run .On the generated polarities aspect classification is performed using Naive bayes, SVM, Decision Tree and KNN.

## IV. RESULTS AND DISCUSSION

Our proposed aspect based sentiment analysis models(Model1 using POS tagging and model2 using TFIDF) have considered reviews on ipod, collected from NLTK python library containing nearly about 800 reviews. These reviews were preprocessed to remove unwanted and noisy data. After preprocessing, aspect extraction, identification of interesting aspects, and separation into positive, negative and neutral will be performed. Accuracy, Precision, Recall, F-Score have been calculated for both the models using Naïve Bayes, SVM, KNN, Decision Tree classifiers and compared with each other. The performance metrics indicate that aspect based sentiment analysis model2 using TFIDF gives higher accuracy on the given input IPod data set.

**The following are the most interesting positive and negative aspects using Model1**

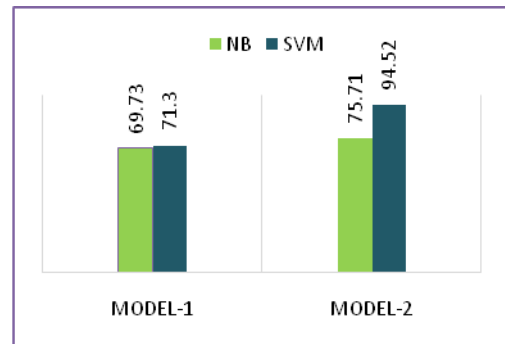| Measure | Naive Bayes | SVM |
|---|---|---|
| Precision | 0.48 | 0.51 |
| Recall | 0.69 | 0.71 |
| F1-score | 0.57 | 0.59 |

**Fig I Performance Metrics**

**The following are the most interesting positive and negative aspects using Model2**

| Measure | Naïve Bayes | K Nearest Neighbor | Decision Tree | SVM |
|---|---|---|---|---|
| Precision | 100 | 80.0 | 85.45 | 92.45 |
| Recall | 61.36 | 100 | 100 | 100 |
| F-Score | 76.05 | 88.88 | 92.15 | 96.07 |
| Accuracy | 75.71 | 80.0 | 87.30 | 94.52 |

**Fig II Performance Metrics**

**Comparison of results**

*Retrieval Number F7935088619/2019©BEIESP*
*DOI: 10.35940/ijeat.F7935.088619*
*Journal Website: www.ijeat.org*

1962

*Published By:*
*Blue Eyes Intelligence Engineering*
*& Sciences Publication*

## V.  CONCLUSIONS

The proposed system extracts the aspects from the product reviews. The reviews are on IPOD. These reviews contain noisy and some unwanted data. Those unwanted data is removed during preprocessing so that training the classifier will be easy. Nouns are considered as the aspect terms and extracted based on POS tagging(Model1) and TFIDF(Model2) By using the threshold value, classification of the aspects into positive, negative and neutral is done. On the extracted aspects Naïve bayes, Support Vector Machine are applied. These classifiers are applied to find out the correct classification of aspects. Finally experimental results are presented. The performance metrics indicate that both Navie Bayes and SVM work similarly on the given input IPod data set.  In future, this model can be implemented using the Fuzzy concepts.

## REFERENCES

1. Nurulhuda Zainuddin, Ali Selamat, and Roliana Ibrahim "Improving Twitter aspect based sentimental analysis using hybrid approach "on 2016.
2. Edison Marrese-Taylor, Juan D. Velasquez, Felipe Bravo-Marquez and Yutaka Matsuo "Identifying Customer Preferences about Tourism Products using an Aspect-Based Opinion Mining Approach" 17th International Conference in Knowledge Based and Intelligent Information and Engineering Systems - KES2013
3. Jingbo Zhu, Huizhen Wang, Muhua Zhu, Benjamin K. Tsou, Matthew Ma "Aspect-based opinion polling from customer reviews" IEEE Transactions on affective computing, Vol. 2, No.1, January-March 2011.
4. D V Nagarjuna Devi, Chinta Kishore Kumar, Siriki Prasad, "A Feature Based Approach for Sentiment Analysis by Using Support Vector Machine" 2016 IEEE 6th International Conference on Advanced Computing.
5. Rajalaxmi Hegde, Dr. Seema. S "Aspect Based Feature Extraction and Sentiment Classification of Review Data sets using Incremental Machine learning Algorithm" 3rd International Conference on Advances in Electrical, Electronics, Information, Communication and Bio-Informatics (AEEICB17)
6. Shwetha Rana, Archana Singh "Comparative Analysis of Sentiment Orientation Using SVM and Naive Bayes Techniques" 2016 2nd International Conference on Next Generation Computing Technologies (NGCT-2016) Dehradun, India 14-16 October 2016
7. Toqir A. Rana and Yu-N Cheah "Improving Aspect Extraction Using Aspect Frequency and Semantic Similarity-Based Approach for Aspect-Based Sentiment Analysis" Springer International Publishing AG 2018 P. Meesad et al. (eds.), Recent Advances in Information and Communication Technology 2017, Advances in Intelligent Systems and Computing 566.
8. Raisa Varghese and Jayasree M "Aspect Based Sentiment Analysis using Support Vector Machine Classifier" 2013 International Conference on Advances in Computing, Communications and Informatics (ICACCI)
9. Kim Schouten, Onne van der Weijde, Flavius Frasincar, and Rommert Dekker "Supervised and Unsupervised Aspect Category Detection for Sentiment Analysis with Co-occurrence Data" IEEE Transactions on cybernetics.
10. Povoda, Lukas, Radim Burget, and Malay Kishore Dutta, "Sentiment analysis based on Support Vector Machine and Big Data", Telecommunications and Signal Processing (TSP), 2016 39th International Conference on, IEEE, 2016.