

Exploring different approaches for sentiment analysis.

Sridhar Ashwin
School of Computing
Dublin City University
Dublin, Ireland
ashwin.sridhar3@mail.dcu.ie

Nair Rohit
School of Computing
Dublin City University
Dublin, Ireland
rohit.nair2@mail.dcu.ie

Foster Jennifer
School of Computing
Dublin City University
Dublin, Ireland
jennifer.foster@dcu.ie

Abstract—*****After the results and analysis*****

Keywords—*****ADD LATER*****

I. INTRODUCTION

The digital era has given users access to the internet and this generates a huge amount of data each second. Every click on the internet generates data which can be captured by the website or the application. In our digital world, data is the new oil. Twitter, Facebook, LinkedIn and other such social media platforms give users a stage to voice their opinion.

People nowadays are keen to know others' opinion about a product/service before they are buying and make their decisions based on the reviews of others. Big E-commerce companies like Amazon or Ebay receive extremely high number of reviews on their products. Several reviews for a specific product, brand, individual personality etc. are very helpful in deciding the perception of people towards the product. Hence it becomes necessary to create algorithms to automate the classification of distinctive reviews on the basis of their polarities particularly: Positive, Negative and Neutral. This automated classification mechanism is referred to as Sentiment Analysis.

Sentiment analysis involves analysing the data and extracting the sentiment attached to it. Natural Language Processing (NLP) enables us to extract emotions from human language. It involves tokenization of sentences, i.e. breaking down sentences into words and deriving meaning from it based on the sentiment value associated with each word.

Each word has a sentiment value attached to it based on its intensity and level of emotions that it conveys. This project explores lexicon based and machine learning based algorithms to compare their accuracy and performance on different data sources.

II. LITERATURE REVIEW

There are systems that have currently implemented both, lexicon and machine learning based algorithms. The most common approach is a mathematical method to derive sentiment from data collected.

Reference [7] extracted tweets from twitter to determine the effects of an average person's tweets over fluctuation of stock prices of a multinational company. Work in this paper focuses on unsupervised learning using sentiment-lexicon based approach. Sentence level approach is being used here. Here, extraction of tweets is done from twitter4j library. The extracted tweets are then pre-processed and cleaned by removing unwanted words in the sentences. The tweets are then split into tokens which are compared with a dictionary of words and scores are assigned to words where a positive word is 1, negative word is -1 and neutral word is 0. The final

sentiment of the tweet is calculated by adding all the scores of the individual tokens.

Xing Fang & Justin Zhan [6] use product reviews from E-commerce websites like Amazon.com to perform both sentence-level and review level categorization. This paper tackles a fundamental problem of sentiment analysis, namely sentiment polarity categorization. They compared Naïve Bayesian, Random Forest and SVM classification. They also implemented parts of speech (POS) tagging to remove parts of speech that do not contribute to the sentiment. The dataset consists of 5.1 million reviews from amazon belonging to 4 major categories: beauty, books, electronics and home. Feature vector formation method is finally used to compute the sentiment of the sentence.

Yige Wang [2] combined POS tagging with Naïve Bayes classifier to train the model which is used to predict the sentiment of movie reviews. With POS tagging, the complexity was reduced by a huge margin as the number of words that had to be trained decreased, but the overall accuracy of the model decreased by 7%. The experiment shows that the sentiment of a sentence or a document is mainly dependent on adjectives and verbs and hence only certain types of words are selected to train the model using POS tagging thereby reducing computational complexity.

Reference [3] compares Naïve Bayes approach and SVM to predict sentiment using review dataset. The dataset they considered of airline reviews which was collected from twitter. They found the accuracy of SVM to be higher than Naïve Bayes, the accuracy of SVM was found to be 83% and Naïve Bayes provided an accuracy of 77%. The result shows that in the case of airline reviews, SVM algorithm gives better results than Naive Bayes algorithm.

H. Parveen and S. Pandey [4] use Naïve Bayes algorithm to reduce the overhead. The current systems lack the capability to handle huge data sets and are very time consuming. The study proposes the use of Hadoop and MapReduce architecture to handle big data. They also used emoticons to extract sentiment from the data. The whole emoticon was converted into its equivalent word and the sentiment associated with the word was used.

U. Kumari, A. Sharma and D. Soni [5] use SVM classification technique to find the sentiment from a product review dataset. The accuracy of the model was tested with datasets of 4 different products. The highest and the lowest accuracy the model achieved were 90.99% and 88.03% respectively. Though the accuracy of the model is high, the datasets used for testing were considerably small.

M. Wongkar and A. Angdresey [1] also use tweets to assess the political sentiment of Indonesian people during the 2019 presidential elections. A comparison was carried out using Naive bayes method, Support-vector networks and K-

Nearest neighbour methods. The method used to collect data from twitter is data crawler in this study. Probability of each word being positive or negative is calculated and multiplied. The result is used to classify the sentence as positive or negative. Twitter is a platform where people mostly use text of less than 140 words to publish their opinion. Hence tweets are one of the most effective datasets for sentiment analysis. In this study, it was found that the Naive Bayes method is the most accurate with 80.1% accuracy.

Bo Pang and Lillian Lee [8] have used IMDB movie review data set. They have used three machine learning methods namely Naïve Bayes, Maximum Entropy classification and support vector machines. The paper classifies movie reviews as either Positive or Negative. This was one of the first attempts to take this approach. Bag of words method was used and Unigrams were taken as features for this classification problem and the end result showed that the system performed well with either of the approaches.

Sida Wang and Christopher D. Manning [9] conclude variants of Naïve Bayes and Support Vector Machines (SVM) behave differently depending on the model variant, feature used and task/dataset. They have identified simple Naïve Bayes and SVM variants that provide great results when compared to most published results on sentiment analysis. A new variant NBSVM that was identified in this paper works well on snippets and longer documents, for sentiment, topic and subject classification and is found to be better than most published results. The work introduced different variants like MNB (Multinomial Naïve Bayes), Multivariate Bernoulli NB (BNB) and NBSVM (Naïve Bayes Support Vector Machines). MNB is found to perform better than BNB in most cases.

III. PROPOSED WORK

This project would explore how lexicon based and supervised machine learning approaches perform for different data sets. Naïve Bayes algorithm will be used for the machine learning based approach. For the lexicon-based approach, a dictionary will be used to identify positive and negative words. The sentiment score will be calculated in the following manner [7]:

$$\frac{\text{Positive Words} - \text{Negative Words}}{\text{Total Number Of Words}}$$

Sentiments from emoticons will also be extracted by the model [4]. Not all words contribute towards sentiments, the data would have to be cleaned before the model begins analysing it. Graphical representation of the analysed data would help in conveying the message in a well ordered and organised manner.

In supervised learning the dataset used to train the model plays a vital role in determining the accuracy of the model. How different datasets affect the accuracy of the system and which approach performs better for each type of data set will be analysed. The accuracy of the model and how it is affected when data sets of different class get introduced will be examined.

Three standard data sets will be used to train and test the models. Stanford Artificial Intelligence Laboratory contains labelled data sets of IMDB movie reviews and amazon user reviews. There are 12000 positive and negative movie reviews for training as well as testing. Amazon data set has

approximately 120000 records and the ratings are on a scale of 5 which was acquired from Stanford Network Analysis Platform. This will have to be converted to positive or negative to train the model. The twitter data set consists of 1.5 million tweets labelled with 1 and 0. 1 being positive and 0 being negative. This data set was acquired from Thinknook.com.

Spyder and Jupyter Notebook will be used for building the model. Additional libraries such as NLTK and Numpy will also be used for analytics and Matplotlib or ggplot can be used for visualisation. Different models will be trained using various algorithms on Python. Implementation of Naïve Bayes and one lexicon-based algorithm will be done.

The template is used to format your paper and style the text. All margins, column widths, line spaces, and text fonts are prescribed; please do not alter them. You may note peculiarities. For example, the head margin in this template measures proportionately more than is customary. This measurement and others are deliberate, using specifications that anticipate your paper as one part of the entire proceedings, and not as an independent document. Please do not revise any of the current designations.

IV. CONCLUSION

Data is being generated by everyone these days. Users post their views and form opinions based on what others share online. The goal of this project is capitalising on this huge source of data, find which approach is better for each dataset and explore the reason behind an approach outperforming the other. Businesses need to know what their customers think about their products and it is humanly impossible to go through each and every review. The goal is to automate the process and give the businesses an overview of the sentiment of the customers towards their product. A decision, if not backed by data, will not yield good results. The results will be evaluated by calculating the accuracy percentage of the model based on the difference in predicted sentiment and actual sentiment. This will help us determine which approach would be better for extracting user sentiment for each subdomain.

REFERENCES

- [1] M. Wongkar and A. Angdresey, "Sentiment Analysis Using Naive Bayes Algorithm Of The Data Crawler: Twitter," in *Fourth International Conference on Informatics and Computing (ICIC)*, Semarang, 2019.
- [2] Y. Wang, "Advanced Naïve Bayes Algorithm Design with Part-of-Speech Tagger on Sentiment Analysis," in *International Conference on Computer Systems, Electronics and Control (ICCSEC)*, Dalian, 2017.
- [3] A. M. Rahat, A. Kahir and A. K. M. Masum, "Comparison of Naive Bayes and SVM Algorithm based on Sentiment Analysis Using Review Dataset," in *8th International Conference System Modeling and Advancement in Research Trends (SMART)*, Moradabad, 2019.
- [4] H. Parveen and S. Pandey, "Sentiment analysis on Twitter Data-set using Naive Bayes algorithm," in *2nd International Conference on Applied and Theoretical*

- [5] U. Kumari, A. Sharma and D. Soni, "Sentiment analysis of smart phone product review using SVM classification technique," in *International Conference on Energy, Communication, Data Analytics and Soft Computing (ICECDS)*, Chennai, 2017.
- [6] X. Fang and J. Zhan, "Sentiment analysis using product review data," *Journal of Big Data*, vol. 2, 2015.
- [7] R. K. Bakshi, N. Kaur, R. Kaur and G. Kaur, "Opinion mining and sentiment analysis," in *3rd International Conference on Computing for Sustainable Global Development (INDIACom)*, New Delhi, 2016.
- [8] B. Pang, L. Lee and S. Vaithyanathan, "Thumbs up? Sentiment Classification using Machine Learning," in *Empirical Methods in Natural Language Processing (EMNLP)*, Philadelphia, 2002.
- [9] S. Wang and C. D. Manning, "Baselines and Bigrams: Simple, Good Sentiment and Topic Classification".
- [10] B. Pang and L. Lee, *Opinion mining and sentiment analysis*, now Publishers Inc., 2008.