

A Mathematical Prediction Model on Sentiment analysis of Twitter data

¹Sudheer Kumar Singh, ²Dr. Prabhat Verma, ³Dr. Pankaj Kumar

Abstract

Sentiment analysis is a process of analyzing textual data to extract sentiment or feeling of textual data on social networks shares by the millions of online people or user's in the Internet world. The current scenario of research in the pasture of sentiment analysis is focused on social network platforms like Twitter, Facebook, etc. The sentiment of data is based on the textual data on social networks. Sentiment analysis is a procedure to extract the positive, negative or neutral sentiment in the form of numeric values range from -1 to 1. It is a sentiment of a user's expression in the form of text about any people or product review, etc. The sentiment of textual data is based on the feature extraction of a text or content. Most of the researchers working on sentiment analysis focus on varieties of features of tweets. In feature extraction, we focused on the combined length of token and character in the textual data. In this paper, we had collected approximately three thousand two hundred tweets from Twitter and extract the sentiment of tweets and count the number of characters or the length of tweets using machine learning concepts in python platforms. We proposed a mathematical model to predict a relationship among the lengths of textual data or the number of characters in tweets of users and its corresponding sentiments using machine learning and linear regression.

Keywords: *sentiment analysis, feature extraction, tokens, social networks, machine learning, Linear regression*

I. INTRODUCTION

social networks have now become a vast field, in each year's millions of peoples are joining the social network sites and communicate with other peoples by sharing their data, information, thoughts, images, etc. Peoples in social network sites use regularly on priority bases due to demand on business, political gathering, news, social media, blogs. A huge amount of data per day circulate on social networks site like Facebook, Instagram, Twitter, etc. the millions of data, information shared on Twitter, data scientist uses these data to their research analyzing these data in different aspects and extract most valuable theories for different areas respectively like, business, politics, movies reassess, product review, etc. Sentiment analysis is one of the superlative fields that are used to analyze data of social networks [3].

¹ Ph.D. Scholar, Dr. A.P.J. Abdul Kalam Technical University, Lucknow, Uttar Pradesh, India.

² Department of Computer Science and Engineering, Harcourt Butler Technological University, Kanpur, India

³ Department of Computer Science and Engineering at Shri Ramswaroop Group of Professional College, Lucknow

The whole world is changing as demand increasing under the present research. The Internet has essentially a requirement of users for a social gathering in the field of social networks. With the exponential growth in social network applications, people are using these platforms to share their data concerning daily issues. Collecting and analyzing peoples' thoughts toward b movies rating, a Product, public services, and so on are vital. Sentiment analysis is a common dialogue preparing task that aims to discover the sentiments behind opinions in texts on varying subjects. In recent years, researchers within the field of sentiment analysis are concerned with analyzing opinions on different topics like movies, commercial products, and daily societal issues. Twitter is an enormously popular micro blog on which clients may voice their opinions [2].Opinion investigation of Twitter data may be a field that has been given much attention over the last decade and involves dissecting tweets (comments.) and the content of those expressions [1]. To learn the method for the execution of sentiment analysis, the challenges that accompany it should be well-known. the essential challenges within the method are text mining, in opinion texts, lexical content alone will be misleading and tough to investigate in general [18].

This paper has dived into sections: Section 2: explores current research work related to sentiment analysis on twitter data. Section 3: Collections of data set from twitter on different resources and implement using this data for sentiment calculation in python platform. Section 4 focuses on the proposed model using linear regression. Section -5 result-oriented graphs and prediction of our proposed method result. Section 6: Conclusion of this paper is based on the result and limitation of datasets. Section 7: future scope and direction of research related to the developed model based on parameter extracted in tweets and reply on any subject or user's prediction on any particular topics.

II. REVIEW OF PREVIOUS RESEARCH WORK

2.1 Social Network Analysis

The social network is a platform that uses for the user's interactions or communications with a totally different view in the current scenario. This type of interaction among online users to share their ideas and feeling in terms of text got a lot of attention. Millions of individuals offer their opinion on dissimilar subjects regularly on social media like Facebook, Twitter, and Instagram It has a number of other applications in different filed on social networks for analyzing data from the web that are used for the different types of businesses. [4]. social network platforms like Twitter could be a rich source of concerning people's opinions and extract their opinions for sentiment analysis [5]. For each tweet, it's necessary to find its sentiment and see the sentiment are positive, neutral or negative. Another challenge with twitter is barely a hundred and forty Characters square measure the limitation of every tweet that causes individuals to use phrases and words that don't seem to be in the language process. As of late Twitter has stretched out the content confinements to 280 characters for each tweet.

2.2 Sentiment Analysis on Twitter

Twitter might be a rich stage to discover concerning People's assessments and feeling identifying with the equivalent or various points as they will convey and impart their insight effectively on the informal organization just as Face book and Twitter. There are comparable or various points of assessment arranged application frameworks that plan to remove individuals' suppositions identifying

with various or similar themes. Twitter is a small-scale blog that processes to generate tweets by the users that are interacted with their users or clients or other systems [17]. Twitter has in excess of 313 million unique customers. It is adding thousands of users per day [18]. The sentiment-based systems have several applications from business to social sciences [5]. In as of now, informal communities, especially Twitter, contains small messages and others could utilize entirely unexpected words and truncations utilized on interpersonal organizations [6]. It is trying to remove their supposition by momentum process frameworks essentially, along these lines a few specialists have utilized deep learning and AI with machine learning system to concentrate and mine the extremity of the content [4]. Sentiment analyses for brief texts like Twitter's posts are difficult due to the length of tweets are restricted in the twitter platform. This is paper we tend to analyze the length of tweets posted by users on twitter.

2.3 Approaches of Sentiment Analysis on Social Network

In the Natural language processing environment (NLP), the sentiment analysis covers miscellaneous phenomena regarding how data about emotion, sentiment, opinion, and social identities is conveyed in language [7]. Many tasks consider a multi-directional model in which sentiments are defined on a polarity basis like positive, negative and neutral. [8]. The foremost application is predicting the polarity of reviews like products, movies, and comments related to any opinion, company, and service reviews [1]. We take the concept of the assumption made in the sentiment analysis on social networks in the previously developed model. But our focus is on the sentiment of user's post or twee user's which has some effect of the sentiment of text in the twitter, and some feature factor of text and their length of tweets [9]. This involves many fields like marketing, political election, launching movies etc. [10]. We have studied the different procedures, techniques. By analyzing the various methods, we identified problem in the field of sentiment analysis that most of the sentiment prediction or calculation based on the features of text [11]. Sentiment analysis is mostly extracted from text using the advanced machine learning method. Machine learning is an application of Artificial Intelligent [24]. In an Artificial intelligent field, Artificial intelligence is a rich field for Machine learning applications [1]. This application is used to classified the unstructured or structured data. Especially in the domain of machine learning, it has many classification method has developed [25]. The method uses different techniques to classify unlabeled data [13]. Classifiers have required test as well as training data. Several examples for classification of the structured or unstructured data using machine learning classifiers are known as maximum Entropy, naive Bayes and support vector machine [12]. Machine-learning is two types first is supervised machine learning and other is the unsupervised learning [25]. Supervised-machine learning methods are required to test as well as training data [14]. It is significant to bring up that training data used in a classifier efficiently will compose very easier to predict any conceptual occurrence in the future [15].

Sentiment sentence extraction and POS tagg.ng.It is recommended that each one objective content should be removed for sentiment analysis, rather than removing objective content [4]. In our study and analysis of various research papers and books, from above study many researchers focused on all subjective content and were extracted sentiment for each text in one sentence [4] For further future analysis, if we have subjective content then we have the sentiment of all sentences by extracting word by word from sentences [4]. Sentiment sentence defines as which contains minimum one negative or positive

word. As formation of sentences, sentences were leading tokenized into separate English words. Each word of a sentence has its grammar role that defines the sentences however; the word is employed [4]. The role of syntactic is also referred to as the components of speech. English language defines the eight components of speech which are referred as the subject, noun, pronoun, adjective word function word, adverb, preposition, conjunction, and also interjection [19]. In the linguistic communication process, part-of-speech (POS) taggers are made-up to classify the words that are supported by their components of part of speech. [20]. As we discussed above the sentiment analysis of a sentence is calculated by some function in which the POS tagging has very significant role[4]. POS tagger is extremely helpful due to the subsequent two reasons: first is: Words like nouns and pronouns typically don't contain any sentiment [22]. The second is: it's able to filter out some of this type of word in the sentences with the use of POS Tagger. From the above discussion, we find a sentiment of words or text using machine learning approach in three phases which are positive, negative or neutral (1,-1,0) [21].but in our proposed method we consider the sentiment of word or text in ranges from [- 1to + 1] using text blob and sentiment method in python programming .

III. DATA SET FOR SENTIMENT EXTRACTION

In this paper, we have collected data from twitter by tweepy API. T this API of Twitter used to data extraction from twitter by utilizing Tweepy provided by twitter itself. The guidelines to extract sentiment from tweets are as follows in the process diagram and using python Programming libraries.

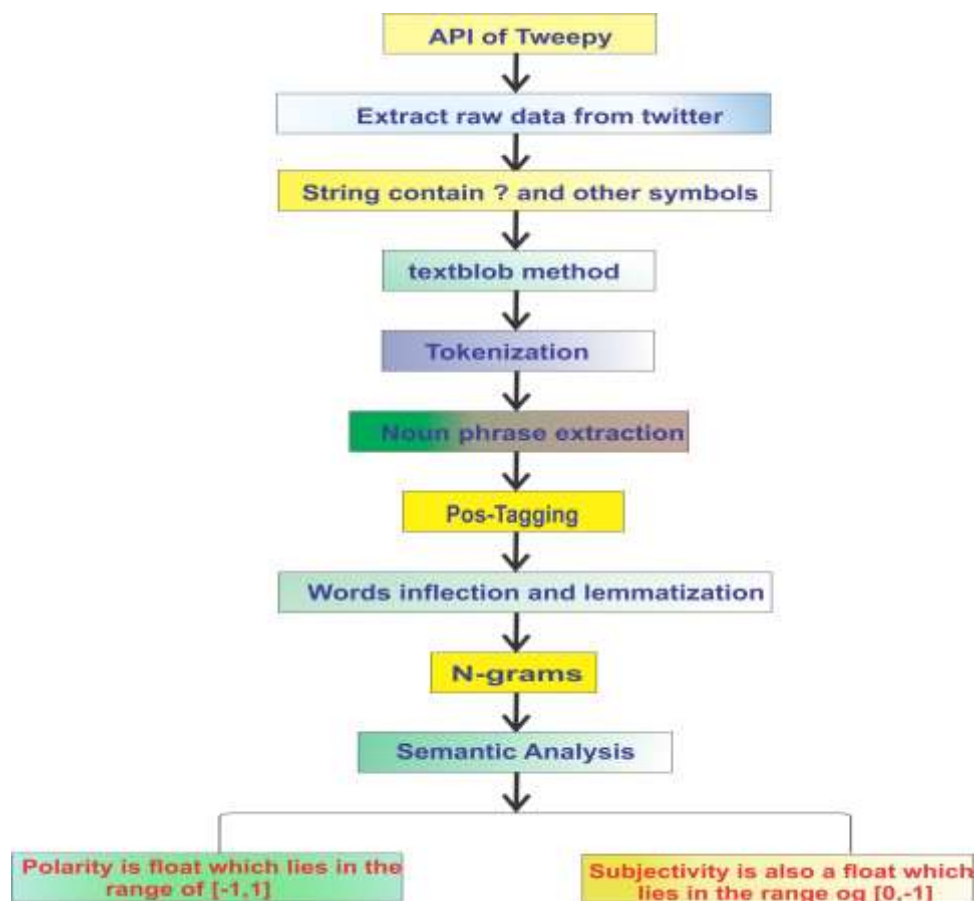


Figure1. Sentiment extraction from twitter dataset

- Collect tweets by downloading the tweets (replies) on followers of any specific user's on twitter.
- Extract data set form the collected data and used it to prepare the training data set.
- Removing the stop words for cleaning the tweets using the text blob method.
- Tokenize each word in the dataset using the regular expression method.
- For each word, text blob and sentiment method calculate the sentiment of text or subjective text according to the rule mentioned in python libraries.

In this paper, we collected tweets (replies) of the followers of a Public figure or Famous personality against the post (tweets) of a Public figure. We collected approximately three thousand two hundred tweets on the twitter platform. All procedure has been implemented using Python language on the Jupiter notebook platform. The Jupyter Notebook is an open-source web application that permits you to make an offer reports that contain live code, conditions, perceptions and account content. Utilizations include information cleaning and change, numerical recreation, factual displaying, information representation, AI, and significantly more. first organized all tweets of followers and Secondly using function extract the length of the tweets (reply). In the table -1, Abbreviation of SOR and LEN are sentiment of reply and length of tweets respectively.

3.1 Data Extracted From Twitter

```
In [18]: In import pandas as pd
data=pd.read_csv("C:/Users/MIRZA/Desktop/Sudheer Sir/work/work/NarendraModi.csv")

In [19]: In data.head()

Out[19]:
      replies  len  tweet
0  Wonderfully expressed by Pranab Dal In/Int is ...  139  Wonderfully expressed by Pranab Dal In/Int is ...
1  RT @IndianNewYork: International Yoga Day at ...  140  Wonderfully expressed by Pranab Dal In/Int is ...
2  RT @Indemb_Muscat: The honourable PM's message...  140  Wonderfully expressed by Pranab Dal In/Int is ...
3  RT @IndianEmbTokyo: Pics of #IDY2019 events or...  139  Wonderfully expressed by Pranab Dal In/Int is ...
4  RT @IndianSwiss: YOGA DAY CELEBRATIONS IN BER...  139  Wonderfully expressed by Pranab Dal In/Int is ...

In [21]: In data.shape,data.dtypes

Out[21]: ((3284, 3), replies  object
len      int64
tweet   object
dtype: object)
```

3.2 Data cleaning and converting into Words by Implementing Textblob and Re Method

```
In [22]: In import textblob
import re

def clean_tweet(tweet):
    """
    Utility function to clean the text in a tweet by removing
    links and special characters using regex.
    """
    return ' '.join(re.sub("(@[A-Za-z0-9]+)|([^0-9A-Za-z \t])|(\#+|\/|\\S+)", "", tweet).split())

def analize_sentiment(tweet):
    """
    Utility function to classify the polarity of a tweet
    using textblob.
    """
    analysis = textblob.TextBlob(clean_tweet(tweet))
    #return str(analysis.sentiment.polarity)
    return(analysis.sentiment.polarity)

In [23]: In l=[]
for tweet in data['replies']:
    l.append(analize_sentiment(tweet))
```

3.3 Implementation using Python Programming on Jupyter Notebook extracts the length and Sentiment of Tweets

```
In [24]: In import numpy as np
data['SOR']=np.array(1)
```

```
In [25]: In lt=[]
for tweet in data['tweet']:
    lt.append(analize_sentiment(tweet))
```

```
In [27]: data.head()

Out[27]:
```

		replies	len	tweet	SOR
0	Wonderfully expressed by Pranab Dal \nInit is ...	139	Wonderfully expressed by Pranab Dal \nInit is ...	0.500000	
1	RT @IndiainNewYork: International Yoga Day at ...	140	Wonderfully expressed by Pranab Dal \nInit is ...	0.242424	
2	RT @Indemb_Muscat: The honourable PM's message...	140	Wonderfully expressed by Pranab Dal \nInit is ...	0.800000	
3	RT @IndianEmbTokyo: Pics of #IDY2019 events or...	139	Wonderfully expressed by Pranab Dal \nInit is ...	0.000000	
4	RT @IndiainSwiss: YOGA DAY CELEBRATIONS IN BER...	139	Wonderfully expressed by Pranab Dal \nInit is ...	0.000000	

EN	40	4	5	3	9	0	39	28	31	36	2	
OR	53	25	7	0	52	53	45	62	26	1	37	34

Table1. Sample of data set having length and sentiment tweets

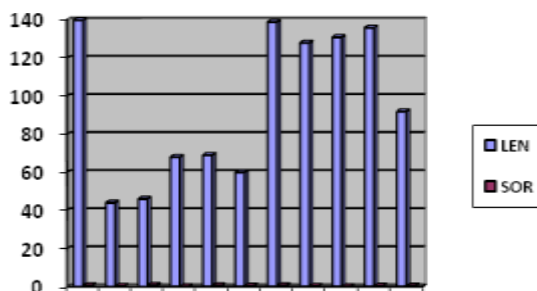


Figure2. Graph between extracted sentiments and length of tweets

IV. PROPOSED MEHTODOLGY

In our proposed model, we use linear regression to predict the relationship between length of tweets of reply and sentiment of tweets.

4.1 Linear Regression

Linear regression is another technique borrowed by machine learning from the field of statistics Regression Analysis is a technique for statistical assessments that empower for the most part of two things:

- **Narrative:** To establish a Relationship between dependent and independent variables can be statistically depicted the strategy for linear regression analysis [23].

• **Opinion:** The estimations of the dependent variables can be approximated from the experiential estimations of the independent variables [23].

4.2 Mathematical Representation of Linear Regression:

Linear regression is a straightforward way to deal with supervised learning. It accepts that the dependence of Y on X_1, X_2, \dots, X_n is linear. The facts confirm that the regression function is never straight. Linear regression always uses a single predictor for establishing a relationship. Let X We accept a model $Y = \beta_0 + \beta_1 X + \varepsilon$ where β_0 and β_1 are two unknown constants that contain the intercepts of y-axis and tangents from the origin and it is also known as coefficients or parameters, and ε is the error notation term, where ε nearly belongs to $N(0, \sigma^2)$.

For the given few estimated values β_0 and β_1 for the proposed model coefficients, we forecast sentiment with the equation $Y = \beta_0 + \beta_1 x$, where y point to a forecast of Y on the origin of $X = x$. It can be shown with some algebra and calculus that this occurs when β_0 and β_1 take the following values:

$$\beta_0 = \bar{Y} - \beta_1 \bar{X}_1 \quad \beta_1 = \frac{\sum(X_i - \bar{X}_1)(Y_i - \bar{Y}_1)}{\sum(X_i - \bar{X}_1)^2}$$

The objective of this research paper is to analysis of data set from this model predict the relationship between the length of tweets and their sentiments [24]. As previous research work has been carried out by using Naive Bayes and Random Forest algorithm to classifying the tweets to estimate the sentiment of the user's tweets on Twitter [23]. For our model we did sentiment analysis using Textblob function in the python platform and make a list of the sentiments for the corresponding user's tweets. Regression method is used to make the prediction model [23]. The results of sentiment analysis are used to predict the relationship between its lengths to sentiments. Our experiment shows to predict the relation by taking data set from previous sentiments of tweets and its corresponding lengths.

In this mathematical model, our assumptions for this model are as described: taking the length of tweets and sentiment of tweets are dependent and independent variables respectively. By using Machine learning algorithms to build a mathematical model, first, we trained this model using training data. Secondly, we apply a linear regression methodology to our text data to analyze data set to predict the desired level of result for future perspectives. In python, we implement our method by calling the linear regression model for the training and test data set. Test whether our result is related to both parameter sets or not. The implementation of linear regression based on the length of tweets and sentiment of tweets is as follows.

```
In [2]: # Importing the dataset and extract dependent and independent variable
dataset = pd.read_csv('Result.csv')
X = dataset.iloc[:, 4:].values
y = dataset.iloc[:, -1].values

In [3]: y

Out[3]: array([139, 140, 140, ..., 73, 137, 140], dtype=int64)
```

Figure3. Importing the dataset and extract dependent and independent variables


```
In [5]: #from sklearn.# Visualising the data set by drawing correlation map
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 1/3, random_state = 0)

In [6]: # Fitting Simple Linear Regression to the Training set
from sklearn.linear_model import LinearRegression
regressor = LinearRegression()
regressor.fit(X_train, y_train)

Out[6]: LinearRegression(copy_X=True, fit_intercept=True, n_jobs=None,
normalize=False)

In [7]: # Predicting the Test set results
y_pred = regressor.predict(X_test)
print(y_pred)

[140. 139. 136. ... 128. 140. 140.]
```

Figure4. Linear regression implementation using correlation map and prediction the test set result.

```
In [8]: # Visualising the Training set results
plt.plot(X_train, y_train, color = 'red')
plt.plot(X_train, regressor.predict(X_train), color = 'blue')
plt.title('Sentiment Analysis of Twitter (Training set)')
plt.xlabel('SOR')
plt.ylabel('LEN')
plt.show()
```

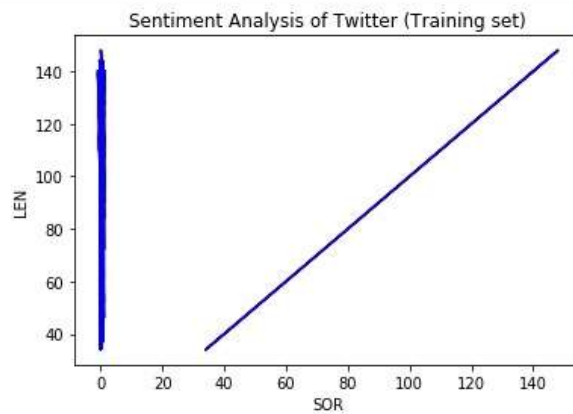


Figure 5. Visualizing the training set results

```
In [10]: #Finding Residuals
from sklearn.metrics import mean_squared_error, mean_absolute_error
print('MAE : ',mean_absolute_error(y_test, y_pred))
print('MSE : ',mean_squared_error(y_test, y_pred))
print('RMSE : ',np.sqrt(mean_squared_error(y_test, y_pred)))

MAE : 1.8555277996581642e-14
MSE : 1.5552200701642587e-27
RMSE : 1.362177594757073e-07
```

figure6. Showing the residuals results as error in calculation

V. CONCLUSION AND FUTURE DIRECTIONS

The objective of this research paper is to represent a review of previous work related to sentiment analysis on the social network using the twitter platform. We focus on Public figure tweets of any topic or issue and their follower's responses (tweets) for a particular subject or issue on twitter and the length Of the tweets, In previous research on sentiment analysis, it is based on sentiments of users(texts). The sentiment of a sentence or content depends on the highlights feature of a content.. Features of text have contained many parameters. As a result features extraction of text depends on different parameters that sentiment of a text depends on different parameters. As we discussed above, we have taken a parameter length extract from the (tweets) replied by the users. We proposed a novel mathematical model using linear regression and machine learning to establish relationships among the length of tweets and their sentiment. Our proposed model implements linear regression using machine learning techniques on the collected dataset over twitter. As a result, our proposed model predicts the relation among the tweets' length and sentiment of respective tweets is not linear, sometimes it is linear and sometimes it behaves like nonlinear. But an average it is represented a nonlinear graph. This model predicts that the sentiment of text depends on many parameters of text; it does not depend on specific one parameter. In the future, we proposed a model, in this model, we find sentiment of tweets before posting it on twitter by checking sentiment has positive or negative. if the sentiment is negative, we have the option to modified the tweets by using keywords having positive sentiments. As we analyze the sentiment of text depends on features of text having emotions or feeling like keywords or words. Using these results in the future, we will develop a game-theoretic model to establish the relationship among sentiments of users on social networks like Twitter using concepts of analysis of the social network.

REFERENCES

- [1] Khan, Abdullah Alsaeedi1 Mohammad Zubair. "A Study on Sentiment Analysis Techniques of Twitter data". *International Journal of Advanced Computer Science and Applications*. Vol. 10, No. 2.(2019)
- [2] Carroll, TarasZagibalov John. "Unsupervised Classification of Sentiment and Objectivity in Chinese Text". *ZagibalovCarroll*. 2008.
- [3] Lee, Sang-Hyun; Lee, Lee-Sac; Hwang, Hyun-Seok." Does Social Opinion Influence Movie Ticket Revenues?: A Case Study" *American Scientific Publishers Advanced Science Letters*, Vol. 23, No. 3, March 2017, pp. 1627-1630.
- [4] Alexander Pak, Patrick Paroubek." Twitter as a Corpus for Sentiment Analysis and Opinion Mining" *European Language Resources Association (ELRA).vol Proceedings of the Seventh International Conference on Language Resources and Evaluation (LREC'10)*, May 2010.
- [5] Soujanya Poria, Erik Cambria, Alexander Gelbukh." Deep Convolutional Neural Network Textual Features and Multiple Kernel Learning for Utterance-Level Multimodal Sentiment Analysis". *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 2539–2544, Lisbon, Portugal, 17-21 September 2015.

- [6] Deepak Kumar, Shivani Aggarwal. "Analysis of Women Safety in Indian Cities Using Machine Learning on Tweets" Amity International Conference on Artificial Intelligence (AICAI) February 2019.
- [7] Pang and Lillian Lee. 2008. Opinion mining and sentiment analysis. *Foundations and Trends in Information Retrieval*, 2(1):1–135.
- [8] Lisa Feldman Barrett and James A. Russell. "Independence and bipolarity in the structure of affect" *Journal of Personality and Social Psychology*, 74(4):967–984, 1998.
- [9] David C. Rubin and Jennifer M. Talerico. "A comparison of dimensional models of emotion". *Memory*, 17(8):802–808(2009).
- [10] Bei Yu, Stefan Kaufmann, and Daniel Diermeier. "Classifying party affiliation from political speech" *Journal of Information Technology and Politics*, 5(1):33–48(2008).
- [11] Muhammad Zubair Asghar¹, Aurangzeb Khan², Shakeel Ahmad¹, Fazal Masud Kundi¹ "A Review of Feature Extraction in Sentiment Analysis" *J. Basic. Appl. Sci. Res.*, 4(3)181-186, (2014).
- [12] K. P. Murphy, "Naive bayes classifiers," University of British Columbia, vol. 18, 2006.
- [13] A. L. Berger, V. J. D. Pietra, and S. A. D. Pietra, "A maximum entropy approach to natural language processing," *Comput. Linguist.*, vol. 22, no. 1, pp. 39-71, 1996.
- [14] A. S. Nugroho, A. B. Witarto, and D. Handoko, "Support vector machine," *Teori dan Aplikasinya dalam Bioinformatika, Ilmu Komputer. com, Indonesia*, 2003.
- [15] Abdullah Alsaeedi¹, Mohammad Zubair Khan² (IJACSA) *International Journal of Advanced Computer Science and Applications*, Vol. 10, No. 2, 2019
- [16] V. Kharde and P. Sonawane, "Sentiment analysis of twitter data: A survey of techniques," *arXiv preprint arXiv:1601.06971*, 2016.
- [17] S. R. Das and M. Y. Chen, "Yahoo! for Amazon: Sentiment extraction from small talk on the web," *Management science*, vol. 53, no. 9, pp. 1375-1388, 2007.
- [18] Manasee Godsay "The Process of Sentiment Analysis: A Study" *International Journal of Computer Applications (0975 – 8887) Volume 126 – No.7, September 2015*.
- [19] Pang B, Lee L. "A sentimental education: Sentiment analysis using subjectivity summarization based on minimum cuts." *Proceedings of the 42Nd Annual Meeting on Association for Computational Linguistics, ACL, '04. Association for Computational Linguistics, Stroudsburg, PA, USA*.
- [20]
- [21] Roth D, Zelenko D, "Part of speech tagging using a network of linear separators". *Coling-Acl, The 17th International Conference on Computational Linguistics*. pp 1136–1142(1998)
- [22] Kristina T "Stanford log-linear part-of-speech tagger". <http://nlp.stanford.edu/software/tagger.shtml>.l(2003)
- [23] Zhang, Lei & Zhao, Liang & Zhang, Xuchao & Kong, Wenmo & Sheng, Zitong & Lu, Chang-Tien. (2018). Situation-Based Interpretable Learning for Personality Prediction in Social Media. 1554-1562. 10.1109/BigData.2018.8622016.

- [24] Yahya Eru Cakra, Bayu Distiawan Trisedya.” Stock Price Prediction using Linear Regression based on Sentiment Analysis”. ICAC SIS 2015/IEEE 978-1-5090-0363-1/15/.
- [25] Norah Fahad Alshammari, Amal Abdullah AlMansour.” State-of-the-art review on Twitter Sentiment Analysis”. IEEE,978-1-7281-0108-8/19/2019 .
- [26] Xing Fang* and Justin Zhan.” Sentiment analysis using product review data” Journal of Big Data 2:5(2015) s