Sentiment Analysis on Twitter Data using Machine Learning Techniques

Muhammad Faizan Siddiqui¹, Faiza Iqbal², Naveed Hussain³

¹H – Cloud & Infrastructure Services, Systems Limited, Lahore, Pakistan
 ²Department of Computer Science, University of Engineering and Technology, Lahore, Pakistan
 ³ Department of Computer Science, University of Central Punjab, Lahore, Pakistan

ABSTRACT

Sentiment analysis is a major field of text mining. It is used to analyze opinions, sentiments, and ideas written in natural language text such as English. This research proposed an application called "Tweemotions" which utilize sentiment analyses to find the opinion about tweets and categorized them as a positive or negative comment. The emotion proposed application used thousands of tweets obtained from NLTK Corpora. Further proposed applications were trained and tested by seven Machine Learning algorithms such as BernoulliNB, MultinomialNB, Logistic Regression, stochastic gradient descent (SGD) classifier, LinearSVC, SVC, and NuSVC. The Tweemotion achieved more than 81% accuracy using MultinomialNB, BernoulliNB, SVC, and NuSVC.

Keywords: Sentiment Analysis, Twitter, Natural Language Processing, Machine Learning, Scikit-learn

Author`s Contribution	Address of Correspondence	Article info.
¹ Data analysis, interpretation, and	Faiza Iqbal	Received: June 17, 2021
manuscript writing, Active participation in	Email: faizaaibal@gmail.com	Accepted: November 13, 2021
data collection and implementation		Published: December 30, 2021
² Conception, synthesis, planning of		
research, ^{3.} Interpretation and discussion		
Cite this article: Siddiqui F. M, Iqbal F, Husso	in N. Sentiment Analysis on Twitter Data	Funding Source: Nil

Cite this article: Siddiqui F. M, Iqbal F, Hussain N. Sentiment Analysis on Twitter Data using Machine Learning Techniques. J. inf. commun. technol. robot. appl.2021; 12(2):30-38

INTRODUCTION

Today, we stand at a time where anyone can make a public statement, not just in a limited geographical area but to the world [1-4]. The internet and social media have advanced at such a high rate that communication of one's opinion can easily be broadcasted with a few clicks on a lightweight device anywhere in the world. There exist many platforms which can deliver these kinds of communications such as Facebook, Instagram, Twitter, etc [5]. Twitter is considered the most authentic platform for delivering social newscasts. People tweet about different topics all the time and their opinions may not matter to all of us but if we look at this from a business point of view, their opinions are what matter the most [6, 7].

Conflict of Interest: Nil

Sentiment Analysis (SA) is the best procedure to discover the sentimental/emotional value of a given segment of natural language [8,9]. This process can be used to classify emotions in three different portions, i.e. Document, Sentence, and Aspect. In the document portion, a complete document is analyzed resulting in a positive or negative sentiment [1,10,11]. In the sentence portion, the same process is applied to find the sentiment of a sentence which can be considered a short document unit. The aspect portion of sentiment analysis work by

finding existing aspects or features of a text and then analyzing the sentiments of all these features. In the age of social media, SA is applied to figure the public opinion on some specific matters so that accurate decisions can be made [12].

The main purpose of this paper is to apply sentiment analysis (SA) to thousands of tweets to know the type of emotional opinion the public holds regarding a certain topic. The goal is to help narrow down the searching efforts of Twitter users and assist them to provide specific results they desire in a compact form. The proposed methodology has adapted Natural Language Processing (NLP) technique and trained the model with the help of Machine Learning algorithms such as BernoulliNB, MultinomialNB, Logistic Regression, and stochastic gradient descent (SGD) classifier, LinearSVC, SVC, and NuSVC. The proposed model is, first, trained and tested using a tweets dataset obtained from NLTK Corpora. It is, then, executed using Twitter data obtained directly from Twitter. The model is trained over a subset of the dataset and afterward, it is tested on the rest of the tweets data to generate the sentimental values.

The remainder of the paper is assembled as follows: Introduction section is followed by section II which elaborated the proposed methodology. In section III, Machine Learning algorithms, used for the analysis, have been described. Section IV presented results in the form of tables and graphical notation. In the last section, V concluded the paper.

METHODOLOGY

In this section, we introduce the existing Scrum model and explain the need for reliability engineering practice.

Sentimental analysis of Twitter data requires preprocessed and cleaned data. The proposed methodology uses Natural Language Processing (NLP) to analyze text with the help of Machine Learning (ML) techniques to train the model [13, 14]. It follows these steps to generate positive or negative results about a particular tweet. Figure 1 demonstrates these steps in flow chart format.

i. Initially, the dataset is in raw format. It is preprocessed and cleaned by removing stop words, tags, URLs, and stem words.

- ii. The next step is feature extraction. In this step, important features from the cleansed data set have been dugout.
- iii. In the next step, the data set is divided into two sets. Based on the existing research [15], we have used 60 percent of the data for training purposes and 40 percent of data for testing purposes.
- iv. Next, machine learning techniques are applied to the dataset to classify positive and negative tweets and analyze the accuracy of the results.
- v. Finally, the accuracy achieved from different machine learning algorithms is compared and graphically projected in the form of a bar graph.
- vi. Figure 2 represents the algorithm of the complete procedure.



Figure. 1. Flow chart of the proposed methodology

Access to twitter's data

We initially tried training our application model with Twitter data that was available in the cleansed format on multiple websites but our primary source was NLTK Corpora [16]. which consisted of 10,020 tweets. 5016 of which were positive and 5004 were negative tweets. But later on, we got access to twitter's data directly from Twitter and the final training and testing were performed on it. Table I represents the stats of the data set.

Table 1: Dataset statistics							
Dataset	Positive Tweets	Negative Tweets	Total				
Training	2988	3024	6012				
Test	2028	1980	4008				

Preprocessing

When data is obtained directly from the source i.e. from Twitter, it is not in an ordered format. Transformation of this unordered data into a well-organized setup is performed in preprocessing step [17, 18]. The data that we have obtained from Twitter is raw. In this form, the analyses performed on such data using machine learning techniques would be extremely difficult and inaccurate [19]. So, we have made a few changes in it to make it easier for the model to adapt to [20]. Hence, our preprocessing consists of the following corrections:

Algorithm: Extract Twitter Sentiment

Begin

Input Tweets where% user_id, hashtag word, @username Term Frequency Inverse Document Frequency =TF IDF MultinomialNB, BernoulliNB, Logistic Regression, Stochastic Gradient Descent =SGD, Support Vector Classifier=SVC, LinearSVC and NuSVC. Output: Confusion Matrix, Classification Report, Accuracy and K-fold validation For each tweet, Do: Procedure Pre-processing(tweet): Remove all the private Twitter symbols (#SpecificWord, @username, retweet (RT)) Remove all numbers, symbols and punctuations, Convert into capital letters into small letters Remove URL ("http://url and https://url") Remove non-ASCII characters and numeric numbers. Substituting multiple spaces with a single space. Contractions and translate into appropriate slang Replace positive emoji with smile word and negative emoji with bad word All the stop words are removed Follow the Tokenization Stemming Return corpus of clean tweet End Procedure Procedure Feature Extract (clean tweet): Count vectorizer = TF_IDF Max Features, N-gram, Max document frequency, Min document frequency End Procedure Procedure Machine learning Classifier (clean tweet): BernoulliNB, MultinomialNB, Logistic Regression, SGD, LinearSVC, SVC, NuSVC End Procedure **Fnd Until** End



- Removed web links
- Removed words joined with "@" and "#"
- Removed retweets (RT)

- Removed duplicated tweets
- Removed characters that do not have ASCII
- Removed non-English words
- Removed punctuation marks and other stop-words
- Removed symbols, emojis, numbers, extra spaces, and punctuations
- Converted to lowercase
- Contraction and conversion of slang into proper dialects
- Replace positive emojis with smile words and negative emojis with bad words
- Corrected spellings
- Applied tokenization and stemming techniques

Feature Extraction

If we wish to know what sentiment (positive/negative) a specific line holds, we must have weightage for each word. Different vectorizers can be used for such a task but the one we intend to use is the "Term Frequency Inverse Document Frequency" (TF-IDF) [21-23]. TF-IDF is a vectorizer (weight assigner) and it is a part of Scikit-learn. We also worked on a different vectorizer known as "Bag of Words" but we are giving priority to TF-IDF because its accuracy is higher and is an advanced version of the Bag of Words [24, 25]. Equations 1-3 represent the procedure to calculate TF-IDF.

Term Frequency =
$$\frac{no.of \ occurrence \ of \ a \ words \ in \ a \ document}{no \ of \ words \ in \ the \ document}$$
 (1)

Inverse Document Frequency = $\frac{no. of documents}{no. of documents containing words}$ (2)

TF-IDF=TF (document, word) \times IDF (word) (3)

N-Gram

The N-Gram procedure will be used to find the features of the tweet for "Supervised Machine Learning Algorithms" [26, 27]. In N-gram, there is an arrangement of N number of tokens from the tweets. The number of N values should be 1, 2, 3, and so on but by default assuming the 1 value. When we examine the value of N to be 1 it is known as a unigram, for N=2, it is known as a bigram, and for n=3, it is known as a trigram, and so on [28]. Consider a tweet in this data set "Sorry guys, no new video this week". If we consider N=2 then it will produce "Sorry guys,", "guys, no", "no new", "new video", and "this week". On the other hand, if we consider N=3 then it will produce "Sorry guys, no", "guys, no new", "no new video", "new video this", and" video this week". These bigram and trigram tokens will be analyzed with further processing. Table II represents the parameters and their values used in the feature extraction process.

Table 2: Feature Extraction						
Parameters	Values					
N-gram range	(1,2) and (1,3)					
Maximum features	7750					
Minimum document	2					
frequency						
Maximum document	0.08					
frequency						

MACHINE LEARNING TECHNIQUES

There are many techniques of Machine Learning classifiers. We have used it to identify sentiment analysis on the bundle of Tweets. These machine learning techniques proposes a solution to convert sentiments into classification record [29, 30]. We have used the classifier to build the model of machine learning and we have used a library related to Python known as Scikit-learn, it is a dominant and very useful open-source machine learning package that offers many classification techniques [31]. The algorithm defines how the model will learn and also how accurate are the results of the model. The algorithms we have used in our model are described in the following subsections.

BernoulliNB Algorithm

"BernoulliNB" classifier is best for countable data values such as "MultinomiaINB". BernoulliNB algorithm is mainly considered to work on boolean features or binary values. Naive Bayes was used for a multivariate Bernoulli distribution of a huge dataset. In this manner, every single class requires samples, which have to be signified in boolean values [32]. We have used BernoulliNB to change inputs of data into a binary object.

The finding rule for Bernoulli is based on Equation 4 as follows.

$$P\left(\frac{xi}{y}\right) = P\left(\frac{xi}{y}\right)xi + \left(1 - P\left(\frac{i}{y}\right)\right)(1 - x)$$
(4)

The BernoulliNB performs the probabilistic approach to train the model and classification techniques use the data that is dispersed by Bernoulli distributions; i.e., it is utilizing numerous type features but one is supposed to be a binary variable. Accordingly, these classifications have required multiple records to be represented as binary-valued feature vectors. Whenever given some other sort of dataset, a BernoulliNB sample of data depends on the binary values [33].

MultinomiaINB Algorithm

The "Multinomial Naive Bayes" classifier has been recognized the text analyzed and moved into classification. MultinomialNB is a probabilistic model dependent on the proposition of Bayes. It computes the possibilities of every text fitting to separate classes and chose the best class with the highest possibilities. The word naive introduce from the supposition that the whole features are free-standing in the class. Although such an autonomy assumption is not normally correct, this Machine Learning technique regularly performs well with comparatively less estimating time. Also, it needs limited data for training and it is exceptionally simple to execute. Naive Bayesian classifier frequently outperforms many sophisticated classification techniques [34, 35].

Logistic Regression Algorithm

"Logistic Regression" is a very powerful Machine Learning technique that is used for binary classification problems such as positive or negative, yes or no, etc [36]. Logistic Regression has performed some predicting analysis techniques and depends on a total number of possibilities [37, 38]. This model has used a very complicated cost identification function. This function can be explained as the "Sigmoid function". It is also called a "Logistic Function" rather than a linear function. Equation 5 represents this function.

$$S(z) = P\left(\frac{1}{1+e^{-z}}\right)$$
(5)

S(z) = output between 0 and 1 (probability estimate) z = input to the function (e.g. y=mx + b) e = base of natural log

The hypothesis of Logistic Regression follows the tendency to restrict the cost identify function somewhere in the range of 0 to 1. So, the Linear functions cannot be represented with a value higher than 1 or lower than 0.

This is impossible for the hypothesis of Logistic Regression [39].

Stochastic Gradient Descent

"Stochastic Gradient Descent" (SGD) is a clear but highly productive approach to managing discriminative learning of linear classifiers under curved misfortune capacities. We have made an only pass over the way the whole trained data. It is managed and controls stability, afterward every update, the project can be repeated onto a circle of radius $1/\lambda$. For every test, we examined a couple of various arrangements for learning rates such as "Logistic Regression" and "Support Vector Machine" [40, 41]. Although Stochastic Gradient Descent (SGD) has been a very important part of the Machine Learning family for a very long time, it has come to be a great deal of consideration only as of late with regards to costly scale learning. It is applied to the huge scope of sparse Machine Learning issues which are constantly experienced in SGD Classifiers [42]. For sparse data, it effectively scales to issues with larger than 105 features and larger than 105 training samples.

SGD has some advantages such as ease of use, implementation, and productivity. On the other hand, SGD has some disadvantages such as it requires various hyperparameters and is sensitive to feature scaling [43].

Linear Support Vector Classification

"Linear Support Vector Classifier" is a very famous Machine Learning technique used to identify model detection, regression, and classification problem. We can use linear kernels in LinearSVC. It performs classification by constructing a N_dimensional hyperplane which ideally isolates data into two or more classes. It has an adaptable technique to select the loss functions and find out penalties and make the scale improved to a large number of samples. Linear SVC performs a "one vs all" plot [44].

Support Vector Classification

"Support Vector Machine" (SVM) is a Supervised Machine Learning technique it has been used for regression and classification tasks. SVM is generally used for classification tasks. It has been used for a statistical classification approach that depends on the maximum distance between each data and is separated by a hyperplane. the hyperplane is used to split data into two different zones. It has the highest distance from the nearest data points in the two different classes. So, we have used a hyperplane to manage tasks in the classification [45]. SVM gives precise reliable and exact classification results [46]. Support Vector Classifier is used for multiple classes and support is applied according to a "one vs one" scheme.

Nu Support Vector Classification

"Nu-Support Vector Classification" (Nu-SVC) behaves like a Support Vector Classifier (SVC) however it utilizes parameters and functions to manage and handle all number of supporting vectors and train the number of errors with the newest parameter V released. The upper bound deals with the number of errors and the lower bound deals with the supporting vector by using parameter V \rightarrow (1,0) [47].

RESULT AND DISCUSSION

After taking the dataset from NLTK Corpora, we have applied pre-processing and feature extraction techniques to them. Experiments on these datasets were performed using seven different classifiers. Tables III and IV demonstrate k-fold cross-validation (bigram and trigram respectively) accuracy results using all mentioned ML classifiers.

Table 3: K-Fold validation (bigram)										
Algorithm	Score1	Score2	Score3	Score4	Score5	Score6	Score7	Score8	Score9	Score10
NB	0.84	0.79	0.83	0.82	0.82	0.84	0.83	0.83	0.80	0.83
Multinomial NB	0.83	0.80	0.82	0.82	0.83	0.83	0.85	0.84	0.81	0.83
Logistic Regression	0.83	0.78	0.79	0.80	0.80	0.84	0.81	0.80	0.80	0.81
classifier	0.81	0.80	0.79	0.80	0.80	0.86	0.81	0.81	0.80	0.64

34

Linear SVC	0. 83	0.79	0.80	0.80	0.81	0.86	0.81	0.81	0.80	0.81
SVC	0.84	0.80	0.80	0.80	0.81	0.85	0.82	0.82	0.81	0.82
NuSVC	0.84	0.79	0.81	0.81	0.80	0.86	0.81	0.81	0.81	0.82

The summarized accuracy results of all ML classifiers using trigram and bigram have been listed in table V. Results show that MultinomialNB has performed best than all used classifiers. The reason for achieving the best accuracy using MultinomialNB is that it is based on the Bayes theorem which works by computing probabilities of the occurrence of multiple text/words belonging to some class based on the occurrence of each text/word and selects the highest probability class. Both BernoulliNB and NuSVC achieve the second-best accuracy results.

Table 4: K-Fold validation (Trigram)											
	Algorithms	Score1	Score2	Score3	Score4	Score5	Score6	Score7	Score8	Score9	Score10
Bernoulli	NB	0.82	0.76	0.81	0.81	0.81	0.83	0.83	0.82	0.78	0.82
Multinomial	NB	0.83	0.79	0.82	0.81	0.83	0.83	0.85	0.84	0.81	0.83
Logistic	Regressio	0.83	0.78	0.80	0.80	0.80	0.83	0.80	0.80	0.80	0.80
SGD	classifie	0.83	0.79	0.78	0.81	0.80	0.85	0.81	0.80	0.81	0.65
Linear	SVC	0. 83	0.80	0.80	0.81	0.80	0.85	0.82	0.80	0.80	0.80

SVC	0.84	0.80	0.80	0.80	0.80	0.85	0.82	0.82	0.81	0.81
NuSVC	0.83	0.80	0.81	0.82	0.80	0.85	0.81	0.81	0.81	0.82

Table 5: Accuracy Results							
Algorithms	Accuracy (trigram)	Accuracy (bigram)					
BernoulliNB	0.806	0.814					
Multinomial NB	0.818	0.824					
Logistic Regression	0.795	0.795					
SGD classifier	0.810	0.807					
Linear SVC	0.808	0.806					
SVC	0.808	0.810					
NuSVC	0.811	0.814					

Table VI shows the results of precision, recall, and F1score of each machine learning algorithm using the bigram and trigram approaches. It is visible that Multinomial NB achieved overall better precision, recall, and F1- score.

Table 6: Results evaluation									
Classifier	Precisio	on	Recall	Recall		F1-Score			
	bi- gram	tri- gram	bi- gram	tri- gram	bi- gram	tri- gram			
BernoulliNB	0.82	0.82	0.81	0.81	0.81	0.80			
Multinomial NB	0.83	0.82	0.82	0.82	0.82	0.82			
Logistic Regression	0.82	0.80	0.81	0.80	0.81	0.79			
SGD classifier	0.81	0.81	0.81	0.81	0.81	0.81			
Linear SVC	0.81	0.81	0.81	0.81	0.81	0.81			
SVC	0.82	0.81	0.81	0.81	0.81	0.81			
NuSVC	0.82	0.81	0.81	0.81	0.81	0.81			



Figure 2: Final accuracy of the ML classifiers (bigram)

Figures 2 and 3 show the accuracy results of all classifiers in the form of a bar graph using the bigraph and trigraph approaches respectively. Similarly, Figures 4 and 5 represent the graphical representation of the precision, recall, and F1-score for each machine learning classifier using bigraph and trigraph.



Figure 3: Final accuracy of the ML classifiers (trigram)



Figure 4: Graphical representation of the percentage of the precision, recall, and F1-Score for every machine learning algorithm using a biograph



Figure 5: Graphical representation of the percentage of the precision, recall, and F1-Score for every machine learning algorithm using trigraph.

CONCLUSION AND FUTURE RECOMMENDATIONS

Twitter is providing valuable blogging services which have been used to find what's going on at any place and any moment. In this paper, we are working that machine learning-related techniques that can be used to predict sentiments on Twitter. We implemented seven different machine learning techniques by using the Python Scikitlearn library to get sentiment analysis from a bundle of Tweets data. Research results show that machine learning methods, such as BernoulliNB, MultinomialNB, Logistic Regression, SGD classifier, LinearSVC, SVC, and NuSVC have improved the accuracy of real-world tweets directly obtained from Twitter and the dataset is publicly available from NLTK corpora. In the future, more valuable results can be obtained using further cleanable text.

REFERENCES

- A. H. Alamoodi et al., "Sentiment analysis and its applications in fighting COVID-19 and infectious diseases: A systematic review," Expert Systems with Applications, vol. 167, p. 114155, 2021/04/01/ 2021.
- [2] K. Garcia and L. Berton, "Topic detection and sentiment analysis in Twitter content related to COVID-19 from Brazil and the USA," Applied Soft Computing, vol. 101, p. 107057, 2021/03/01/ 2021.
- [3] Z. Jianqiang and G. Xiaolin, "Comparison Research on Text Preprocessing Methods on Twitter Sentiment Analysis," IEEE Access, vol. 5, pp. 2870-2879, 2017.
- [4] M. Kanakaraj and R. M. R. Guddeti, "Performance analysis of Ensemble methods on Twitter sentiment analysis using NLP techniques," in Proceedings of the 2015 IEEE 9th International Conference on Semantic Computing (IEEE ICSC 2015), 2015, pp. 169-170.
- [5] A. Arora, S. Bansal, C. Kandpal, R. Aswani, and Y. Dwivedi, "Measuring social media influencer index- insights from Facebook, Twitter and Instagram," Journal of Retailing and Consumer Services, vol. 49, pp. 86-101, 2019/07/01/ 2019.
- [6] T. M. Nisar and M. Yeung, "Twitter as a tool for forecasting stock market movements: A short-window event study," The Journal of Finance and Data Science, vol. 4, no. 2, pp. 101-119, 2018/06/01/ 2018.
- [7] D. Valle-Cruz, V. Fernandez-Cortez, A. López-Chau, and R. Sandoval-Almazán, "Does Twitter Affect Stock Market Decisions? Financial Sentiment Analysis During Pandemics: A Comparative Study of the H1N1 and the COVID-19 Periods," Cognitive Computation, vol. 14, no. 1, pp. 372-387, 2022/01/01 2022.
- [8] K. Sailunaz and R. Alhajj, "Emotion and sentiment analysis from Twitter text," Journal of Computational Science, vol. 36, p. 101003, 2019/09/01/ 2019.
- [9] L. Yue, W. Chen, X. Li, W. Zuo, and M. Yin, "A survey of sentiment analysis in social media," Knowledge and Information Systems, vol. 60, no. 2, pp. 617-663, 2019/08/01 2019.

- M. Singh, A. K. Jakhar, and S. Pandey, "Sentiment analysis on the impact of coronavirus in social life using the BERT model," Social Network Analysis and Mining, vol. 11, no. 1, p. 33, 2021/03/19 2021.
- [11] J. Wehrmann, W. Becker, H. E. L. Cagnini, and R. C. Barros, "A character-based convolutional neural network for languageagnostic Twitter sentiment analysis," in 2017 International Joint Conference on Neural Networks (IJCNN), 2017, pp. 2384-2391.
- [12] W. Budiharto and M. Meiliana, "Prediction and analysis of Indonesia Presidential election from Twitter using sentiment analysis," Journal of Big Data, vol. 5, no. 1, p. 51, 2018/12/19 2018.
- [13] V. N. Khuc, C. Shivade, R. Ramnath, and J. Ramanathan, "Towards building large-scale distributed systems for Twitter sentiment analysis," presented at the Proceedings of the 27th Annual ACM Symposium on Applied Computing, Trento, Italy, 2012. [Online]. Available: https://doi.org/10.1145/2245276.2245364.
- [14] M. Rathi, A. Malik, D. Varshney, R. Sharma, and S. Mendiratta, "Sentiment Analysis of Tweets Using Machine Learning Approach," in 2018 Eleventh International Conference on Contemporary Computing (IC3), 2018, pp. 1-3.
- [15] S. Elbagir and J. Yang, "Sentiment Analysis of Twitter Data Using Machine Learning Techniques and Scikit-learn," presented at the Proceedings of the 2018 International Conference on Algorithms, Computing and Artificial Intelligence, Sanya, China, 2018. [Online]. Available: https://doi.org/10.1145/3302425.3302492.
- [16] "Natural Language Toolkit," T. Aarsen, Ed., ed, 2019.
- [17] J. Samuel, G. G. M. N. Ali, M. M. Rahman, E. Esawi, and Y. Samuel, "COVID-19 Public Sentiment Insights and Machine Learning for Tweets Classification," Information, vol. 11, no. 6, p. 314, 2020.
- [18] A. B. Soliman, K. Eissa, and S. R. El-Beltagy, "AraVec: A set of Arabic Word Embedding Models for use in Arabic NLP," Procedia Computer Science, vol. 117, pp. 256-265, 2017/01/01/ 2017.
- [19] B. Samal, A. K. Behera, and M. Panda, "Performance analysis of supervised machine learning techniques for sentiment analysis," in 2017 Third International Conference on Sensing, Signal Processing and Security (ICSSS), 2017, pp. 128-133.
- [20] A. M. Ramadhani and H. S. Goo, "Twitter sentiment analysis using deep learning methods," in 2017 7th International Annual Engineering Seminar (InAES), 2017, pp. 1-4.
- [21] R. Ahuja, A. Chug, S. Kohli, S. Gupta, and P. Ahuja, "The Impact of Features Extraction on the Sentiment Analysis," Procedia Computer Science, vol. 152, pp. 341-348, 2019/01/01/ 2019.
- [22] K. M. Alomari, H. M. ElSherif, and K. Shaalan, "Arabic Tweets Sentimental Analysis Using Machine Learning," Cham, 2017: Springer International Publishing, pp. 602-610.
- [23] F. Rustam, I. Ashraf, A. Mehmood, S. Ullah, and G. S. Choi, "Tweets Classification on the Base of Sentiments for US Airline Companies," Entropy, vol. 21, no. 11, p. 1078, 2019.
- [24] S. Ahmad, M. Z. Asghar, F. M. Alotaibi, and I. Awan, "Detection and classification of social media-based extremist affiliations using sentiment analysis techniques," Human-centric Computing and Information Sciences, vol. 9, no. 1, p. 24, 2019/07/01 2019.
- [25] M. Wang, D. Cao, L. Li, S. Li, and R. Ji, "Microblog Sentiment Analysis Based on Cross-media Bag-of-words Model," presented at the Proceedings of International Conference on Internet Multimedia Computing and Service, Xiamen, China, 2014. [Online]. Available: https://doi.org/10.1145/2632856.2632912.
- [26] M. Ghiassi, J. Skinner, and D. Zimbra, "Twitter brand sentiment analysis: A hybrid system using n-gram analysis and dynamic artificial neural network," Expert Systems with Applications, vol. 40, no. 16, pp. 6266-6282, 2013/11/15/ 2013.

Journal of Information Communication Technologies and Robotic Applications http://www.jictra.com.pk/index.php/jictra, pISSN: 2523-5729, eISSN: 2523-5739

- [27] S. Xiong, H. Lv, W. Zhao, and D. Ji, "Towards Twitter sentiment classification by multi-level sentiment-enriched word embeddings," Neurocomputing, vol. 275, pp. 2459-2466, 2018/01/31/ 2018.
- [28] J. Awwalu, A. A. Bakar, and M. R. Yaakub, "Hybrid N-gram model using Naïve Bayes for classification of political sentiments on Twitter," Neural Computing and Applications, vol. 31, no. 12, pp. 9207-9220, 2019/12/01 2019.
- [29] W. Medhat, A. Hassan, and H. Korashy, "Sentiment analysis algorithms and applications: A survey," Ain Shams Engineering Journal, vol. 5, no. 4, pp. 1093-1113, 2014/12/01/ 2014.
- [30] J. Singh, G. Singh, and R. Singh, "Optimization of sentiment analysis using machine learning classifiers," Human-centric Computing and Information Sciences, vol. 7, no. 1, p. 32, 2017/12/11 2017.
- [31] R. Wagh and P. Pune, "Survey on Sentiment Analysis using Twitter Dataset," in 2018 Second International Conference on Electronics, Communication, and Aerospace Technology (ICECA), 2018, pp. 208-211.
- [32] X. Jia, B. Hu, B. P. Marchant, L. Zhou, Z. Shi, and Y. Zhu, "A methodological framework for identifying potential sources of soil heavy metal pollution based on machine learning: A case study in the Yangtze Delta, China," Environmental Pollution, vol. 250, pp. 601-609, 2019/07/01/ 2019.
- [33] G. Lombardo, P. Fornacciari, M. Mordonini, M. Tomaiuolo, and A. Poggi, "A Multi-Agent Architecture for Data Analysis," Future Internet, vol. 11, no. 2, p. 49, 2019.
- [34] F. Ali et al., "Transportation sentiment analysis using word embedding and ontology-based topic modeling," Knowledge-Based Systems, vol. 174, pp. 27-42, 2019/06/15/ 2019.
- [35] I. Segura-Bedmar, C. Colón-Ruíz, M. Tejedor-Alonso, and M. Moro-Moro, "Predicting of anaphylaxis in big data EMR by exploring machine learning approaches," (in eng), Journal of biomedical informatics, vol. 87, pp. 50-59, Nov 2018.
- [36] N. F. F. da Silva, E. R. Hruschka, and E. R. Hruschka, "Tweet sentiment analysis with classifier ensembles," Decision Support Systems, vol. 66, pp. 170-179, 2014/10/01/ 2014.
- [37] W. P. Ramadhan, S. T. M. T. A. Novianty, and S. T. M. T. C. Setianingsih, "Sentiment analysis using multinomial logistic regression," in 2017 International Conference on Control, Electronics, Renewable Energy and Communications (ICCREC), 2017, pp. 46-49.
- [38] S. Symeonidis, D. Effrosynidis, and A. Arampatzis, "A comparative evaluation of pre-processing techniques and their interactions for Twitter sentiment analysis," Expert Systems with Applications, vol. 110, pp. 298-310, 2018/11/15/ 2018.
- [39] W. Cheng and E. Hüllermeier, "Combining instance-based learning and logistic regression for multilabel classification," Machine Learning, vol. 76, no. 2, pp. 211-225, 2009/09/01 2009.
- [40] A. M. Alayba, V. Palade, M. England, and R. Iqbal, "Arabic language sentiment analysis on health services," in 2017 1st International Workshop on Arabic Script Analysis and Recognition (ASAR), 2017, pp. 114-118.
- [41] D. Effrosynidis, S. Symeonidis, and A. Arampatzis, "A Comparison of Pre-processing Techniques for Twitter Sentiment Analysis," Cham, 2017: Springer International Publishing, pp. 394-406.
- [42] V. Patlolla. (2017). How to make SGD Classifier perform as well as Logistic Regression using Parfit [Online]. Available: https://towardsdatascience.com/how-to-make-sgd-classifierperform-as-well-as-logistic-regression-using-parfit-cc10bca2d3c4.
- [43] B. Zadrozny and C. Elkan, "Transforming classifier scores into accurate multiclass probability estimates," presented at the Proceedings of the Eighth ACM SIGKDD international conference

on Knowledge discovery and data mining, Edmonton, Alberta, Canada, 2002. [Online]. Available: https://doi.org/10.1145/775047.775151.

- [44] A. Hasan, S. Moin, A. Karim, and S. Shamshirband, "Machine Learning-Based Sentiment Analysis for Twitter Accounts," Mathematical and Computational Applications, vol. 23, no. 1, p. 11, 2018.
- [45] Y. Al Amrani, M. Lazaar, and K. E. El Kadiri, "Random Forest and Support Vector Machine based Hybrid Approach to Sentiment Analysis," Procedia Computer Science, vol. 127, pp. 511-520, 2018/01/01/ 2018.
- [46] H.-J. Cho and M.-T. Tseng, "A support vector machine approach to CMOS-based radar signal processing for vehicle classification and speed estimation," Mathematical and Computer Modelling, vol. 58, no. 1, pp. 438-448, 2013/07/01/ 2013.
- [47] O. Abdelwahab, M. Bahgat, C. J. Lowrance, and A. Elmaghraby, "Effect of training set size on SVM and Naive Bayes for Twitter sentiment analysis," in 2015 IEEE International Symposium on Signal Processing and Information Technology (ISSPIT), 2015, pp. 46-51.

38