**Special Issue - 2022**

**International Journal of Engineering Research & Technology (IJERT)**
**ISSN: 2278-0181**
**CCICS - 2022 Conference Proceedings**

# An Effective Approach for Mental Health Prediction Using Machine Learning algorithm

T. E. Ramya
Assistant Professor
Department of Computer Science and Engineering
Kongu Engineering College, Erode, Tamil Nadu, India.

S. Sindhupriya
MSc Software Systems
Department of Computer Technology - PG Kongu Engineering College, Erode, Tamil Nadu, India.

*Abstract*—The main objective of mental health prediction using machine learning is to manage and detect the social network mental disorders (SNMDs) based on the twitter data. It is aimed to quantify features and patterns from twitter to know the symptoms and risk factors of mental disorders by using methods of machine learning. In the previous system, the user would be capable of logging the form and fill it which had questions based on the data gathered. It is hard to recognize the social network mental disorders because the mental status cannot be directly examined from the online social activity logs. It leads to produce an inaccurate result. In the proposed system, user comments are gathered from twitter and trained by using a convolution neural network. It could easily predict the features of a particular class and also consumes less time that helps in predicting mental status and emotion behavior of specific users. Mental health prediction is developed by using python where Spyder serves as environmental setups. The model is executed with test data to provide a relatively accurate rate of predicting mental status of the users by their tweets.

*Keywords*—*Health, Stress, Support Vector Machine (SVM), K-Nearest Neighbor (KNN), Convolution Neural Network (CNN-LSTM), tweets.*

## I. INTRODUCTION

Mental Disability enables the result for the shortcomings under the branch of brain chemistry. The interpretation of mental wellbeing is over-demanding to recognize and also to prescribe the therapies to the patients with a differed mental behavior [13]. The maximum number of the individual users is liable to stress, while some of them are pretentious to depression for being different reasons. The board of administration of World Health Organization (WHO) determined the depression [3] [7] is found to be the preeminent source of global disease strain [1]. Mental health prediction takes machine learning algorithms [9] and it compares with their accuracy rate separately with the classification algorithm. The algorithms taken for processing are:

- Support Vector Machine (SVM)

- K-Nearest Neighbor (KNN)

- Convolution Neural Network (CNN-LSTM)

The target is to categorize the emotional behavior from the target set of the individuals. For making this possible, it has been implemented on the tweets that has been tweeted by the user. It mainly concerns the responses given by the individuals and has a deeper understanding in classifying the mental factors [6] [7]. SVM is the supervised algorithm [9] of machine learning, which includes used the classification or regression [4] challenges. However, it is more specifically used in classification problem.

The text analysis method based on CNN can obtain the important features of text through the pooling. It will be mostly beneficial for the users and employers to set up the higher awareness about work related mental illness [8] [10] and to overcome the depression [3]. The proposed methodology is based on KNN classification algorithm, which shows an improvement over one of the existing methodologies which is based on SVM classification algorithm [5].

These algorithms have been executed using a machine learning tool called SPYDER. In previous systems, the execution was done by these algorithms using data mining techniques [5] [16] using WEKA environment [15] and it is stated that the Random Tree has more accuracy rate than the others. The dataset has been extracted from KAGGLE. The main objective of this mental health prediction [14] [12] is to prove the result, by combining these three algorithms under some circumstances using SPYDER tool, which can increase the accuracy rate and time taken to execute.

The new dataset set has been selected and the same procedure is followed in the paper. The accuracy rate of the entire algorithm was noted and also the new classification algorithm's accuracy rate is noted and is compared with the existing algorithms. The accuracy rate and the time taken to execute the algorithms were differed and the found classification algorithms have the better execution time and accuracy rate.

Though the rest of the paper is organized as different sections. Section II defines the proposed algorithm which illustrates the framework for Mental Health Prediction approach in detail. Section III presents the evaluation results of these algorithms, and Section IV enables the concluding remarks.

## II. RELATED WORK

The text analysis method based on CNN can obtain the important features of text through the pooling. It will be mostly beneficial for the users and employers to set up the higher awareness about work related mental illness [8] [10] and to overcome the depression [3].

The target is to categorize the emotional behavior from the target set of the individuals. For making this possible, it has been implemented on the tweets that has been tweeted by the user. It mainly concerns the responses given by the individuals and has a deeper understanding in classifying the mental factors [6] [7].

SVM is the supervised algorithm [9] of machine learning, which includes used the classification or regression [4] challenges. However, it is more specifically used in classification problem.

In previous systems, the execution was done by these algorithms using data mining techniques [5] [16] using WEKA environment [15] and it is stated that the Random Tree has more accuracy rate than the others. The dataset has been extracted from KAGGLE.

The main objective of this mental health prediction [14] [12] is to prove the result, by combining these three algorithms under some circumstances using SPYDER tool, which can increase the accuracy rate and time taken to execute.

## III. PROPOSED ALGORITHM

Data is collected by using a twitter API. Then the collected Kaggle dataset with different labels [13] has to be preprocessed by using preprocessing techniques such as tokenization and stop removal method. Though the system flow diagram for the proposed system is shown in Fig. 1. It involves collection of data, preprocessing, encoding the data, processing using trained datasets, testing the algorithm and predicting the results which shown is system flow diagram. The labeled data is encoded for better prediction.
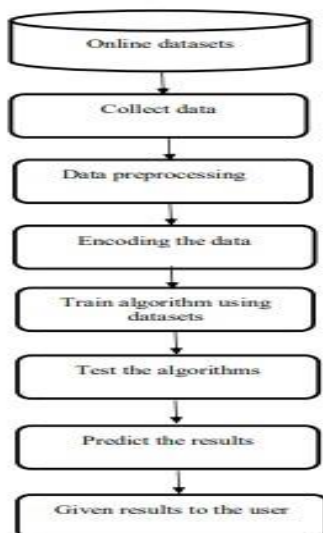


Fig 1 System flow diagram for the proposed system

### A. Data Acquisitions

The dataset has been obtained from the open source Kaggle repository. The fields or labels [13] that taken into major consideration includes the tweets done by the particular individual, date and time on which the tweet has been posted, that specific network, gender and age. The dataset [2] also includes additional information that has been collected for mental behavior for analyzing but not all the parameters are used.

### B. Tweets Preprocessing

Preprocessing is one of the important steps in Data Mining [11] [17]. The tweets after extraction undergo preprocessing to produce the clean dataset. The preprocessing library in python for tweets is preprocessor whose architecture is shown in fig 2. It is very helpful in analyzing mental and emotional behavior of an individual. In preprocessing, there are many steps like stop words removal, punctuations removal and URLs, stemming and finally converted it to a unique format. The library makes efficient cleaning of tweet data, parse them or the tweets can be tokenized. It mainly deals with the following constraints such as URLs, mentions, reserved words, emojis and smileys. The URLs and mentions are first removed furthermore and is followed by removing the digits, punctuations and stop words by using Natural language tool kit (NLTK) [18] which is the best library in processing the twitter data. For accurate information retrieval, the proper dataset with accurate features is obtained
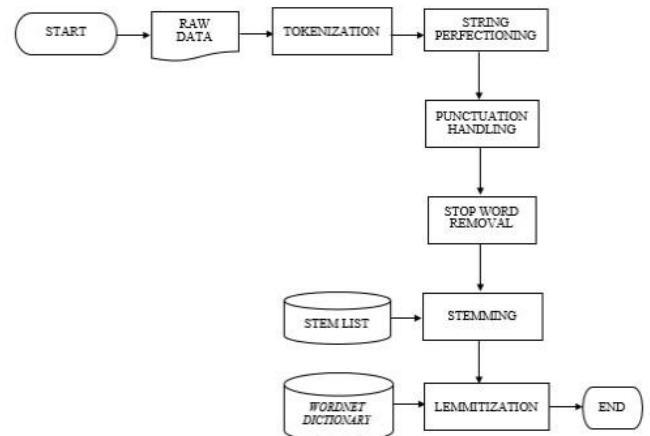


Fig 2 Architecture for preprocessing

### C. Feature extraction

The major aim of feature extraction is to reduce the number of redundant tweets from the preprocessed datasets. The features are used in creating dictionaries which is further assigned with scores in analyzing the peculiar behavior of an individual while tweeting. The variables are created using 0 and 1 association scores which helps in calculating the emotions that contributes to decide mental factor. The features extracted from the dataset have been classified under various machine learning algorithms.

### D. Analysis Of Twitter Data

The label encoded data helps much in the classification process. In SVM [5], the tweets first transformed into vectors. SVM helps in defining the best fit which further classifies which label belong to the concerned category. The main problem involved in text classification is that it includes the mixture of both the characters and words. It is rectified by feeding the numerical representation of the words to make predictions. When CNN is implemented on the cleaned data, the result obtained is detected whether a specialized pattern is

**Special Issue - 2022**

**International Journal of Engineering Research & Technology (IJERT)**
**ISSN: 2278-0181**
**CCICS - 2022 Conference Proceedings**

observed or not. The kernel size is varied and the outputs are finally concatenated. When the length of feature vectors is same, decision trees is involved in categorizing trees with similar patterns.

*E. Visualization*

The accuracy value obtained from each classification process is taken into consideration in choosing the best model for mental behavior analysis. The performance metrics such as precision, recall, F1 score and confusion matrix is evaluated and plotted on the graph for providing better visualization.

$$P = TP\ TP + FP \qquad (1)$$
$$R = TP\ TP + FN \qquad (2)$$
$$F-Measure = \frac{2*\text{Precision}*\text{Recall}}{\text{Precision}+\text{Recall}} \qquad (3)$$

## IV. EXPERIMENT AND RESULT

The experimental setup was completed in Windows, Anaconda environment, python version 3.6.6. The Spyder was used, which enables the ease of application access. The accuracy is calculated for the proposed system which proves that LSTM yields a very good accuracy of 83.76% when compared to the existing preliminary models.

$$A = \frac{TP + TN}{N} \qquad (4)$$

By calculating the sentiment score and the F1 score, the analysis of the user's mental health and emotions are predicted in the form of plot and text outlets. Table-1 shows the accuracy measure of proposed algorithms in comparison with existing algorithms.

TABLE 1 Accuracy measure for existing system and proposed system

| System | Algorithm | Accuracy% | Precision | Time Taken |
|---|---|---|---|---|
| Existing System | Random forest | 82.2 | 0.827 | 0.3 sec |
| | Decision Tree | 79.3 | 0.793 | 0.7 sec |
| | Naïve bayes | 78.7 | 0.787 | 0.5 sec |
| Proposed System | SVM and CNN | 83.76 | 0.837 | 0.2 sec |

Though the analyzation of twitter data for the processing is shown in fig. 3. The LSTM results enable the effective accuracy rate of processing which is shown in fig. 4.
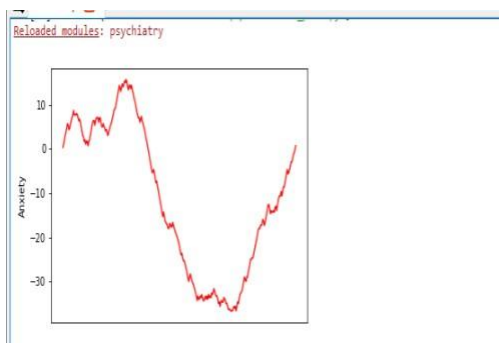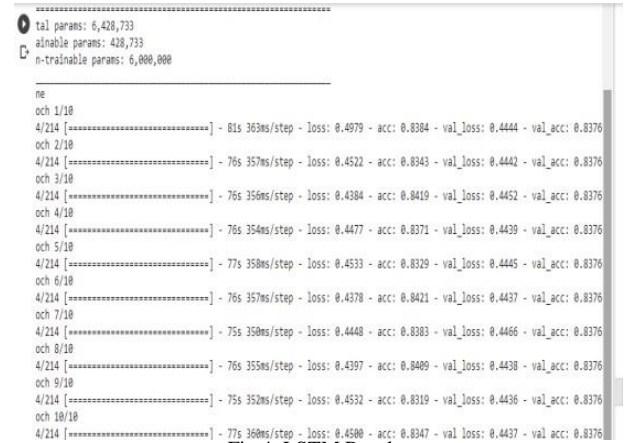


Fig 3. Analyzing Twitter Data



Fig 4. LSTM Results

The frequency chart defines the visualization arrangement of positive and negative commands in access to frequency and word count which is shown in figure 5.The Accuracy rate defines the effective processing rate that LSTM yields the accuracy rate of 83.76% which is found effective in comparison with existing algorithms and is shown in the figure. 6 effectively which is the higher accuracy rate than the existing algorithms.
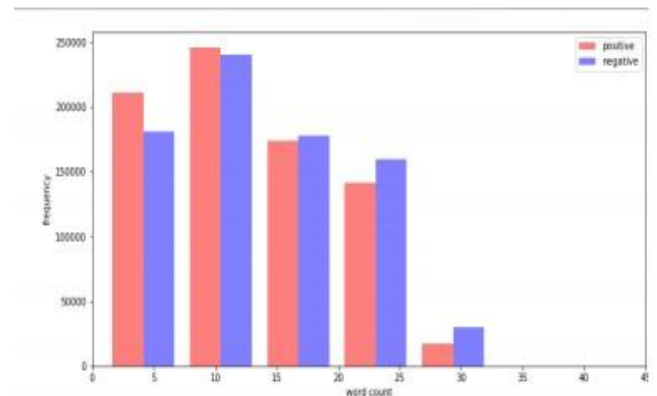


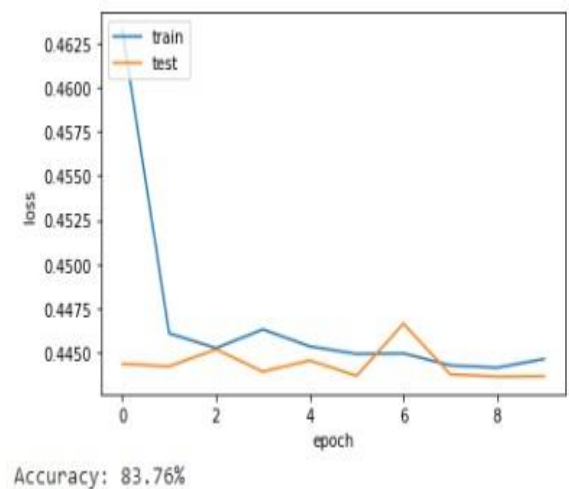Fig 5. Frequency chart for positive and negative commands



Fig 6. Accuracy rate and result analysis

**Special Issue - 2022**

**International Journal of Engineering Research & Technology (IJERT)**
**ISSN: 2278-0181**
**CCICS - 2022 Conference Proceedings**

## V. CONCLUSION

The association score is calculated and has been plotted in the dataset for emotions such as anger, anticipation, fear, joy, and openness. The accuracy value plotted helps in better understanding of model used in classifying the text and also in creating dictionaries. The problem of manual filling questions by the users based on the observed dataset has been overcome.

The suggested method effectively helps in analyzing the emotional and mental behavior of the particular individual from the tweets that has been posted by the user. The models used in the classification of text such as Support Vector Machine, decision tree , Convolutional Neural Network and K Nearest Neighbor effectively classifies the labels that belongs to particular class. The output with the association scores for the corresponding emotions has been calculated with the accuracy of 83.76% and is stored in the new dataset for further analysis of emotional behavior.

## VI. REFERENCES

[1] D.Filip & C. Jesus. (2015). A Neural Network Based Model for Predicting Psychological Conditions International Conference on Brain Informatics and Health 252-261.

[2] Dataset: https://osmihelp.org/research: 2014 Dataset

[3] DEPRESSION:AGlobalCrisis,https://www.who.int/mental_health/manage ment/depression/wfmh_paper_depression_wmhd_2012.pdf March '12

[4] Deziel, M., Olawo, D., Truchon, L., &Golab, L. Analyzing the Mental Health of Engineering Students using Classification and Regression. EDM (2013).

[5] Husseini Orabi, Ahmed & S. G. Alonso, I. Torre-Díez, S. Hamrioui, M.l LópezCoronado, D. C. Barreno, L. M. Nozaleda, and M. Franco. Data Mining Algorithms and Techniques in Mental Health: A Systematic Review. J. Med. Syst. Vol. 42, 9 (September 2018), 1–15

[6] Lumpur, 2018, pp. 1-5. T. Al-Moslmi, N. Omar, S. Abdullah and M. Albared, "Approaches to Cross-Domain Sentiment Analysis: A Systematic Literature Review", IEEE Access, vol. 5, pp. 16173-16192, 2017.

[7] M. A. Haziq Megat S'adan, A. Pampouchidou and F. Meriaudeau, "Deep Learning Techniques for Depression Assessment," 2018 International Conference on Intelligent and Advanced System (ICIAS), Kuala Buddhitha, Prasadith & Husseini Orabi, Mahmoud & Inkpen, Diana. (2018). DeepLearning for Depression Detection of Twitter Users. 88-97. 10.18653/v1/W18-0609.

[8] M. P. Dooshima, E. N. Chidozie, B. J. Ademola, O. O. Sekoni, I. P. Adebayo, A Predictive Modelfor the Risk of Mental Illness in Nigeria Using Data Mining, International Journal of Immunology.Vol. 6, No. 1, 2018, pp. 5-16. 35

[9] M. Srividya, M. Subramaniam and B. Natarajan, "Behavioral Modeling for Mental Health using Machine Learning Algorithms" "Journal of Medical Systems" Vol. 42(5):88 May 2018.

[10] U. S. Reddy, A. V. Thota and A. Dharun, "Machine Learning Techniques for Stress Prediction in Working Employees," 2018 IEEE International Conference on Computational Intelligence and Computing Research (ICCIC), Madurai, India, 2018, pp. 1-4.

[11] Tomar, D., & Agarwal, S. (2013). A survey on Data Mining approaches for Healthcare. International Journal of Bio-Science and Bio- Technology, 5(5), 241-266.

[12] Bhakta, I and Sau, A. (2016). Prediction of Depression among Senior Citizens using Machine Learning Classifiers. International Journal of Computer Applications Vol. 144 No. 7 pp.11–16.

[13] A. R. Subhani, W. Mumtaz, M. N. B. M. Saad, N. Kamel and A. S.Malik, "Machine Learning Framework for the Detection of Mental Stress at Multiple Levels," in *IEEE Access*, vol. 5, pp. 13545-13556,2017.

[14] Sandhya P, M. Kantesaria "Prediction of Mental Disorder for employees in IT Industry", International Journal of Innovative Technology and Exploring Engineering (IJITEE), Volume-8 Issue-6S, April 2019.

[15] WEKA: https://www.cs.waikato.ac.nz/ml/weka/

[16] Han, Kamber, "Data Mining Concepts and Techniques", Morgan Kaufmann 3nd Edition.

[17] P. N. Tan, M. Steinbach, Vipin Kumar, "Introduction to Data Mining", Pearson Education.

[18] Mihalcea R, Tarau P. Text-rank: bringing order into texts. In: Proceeding of the conference on "empirical methods in natural language processing " 2004: 404–411.