

PAPER • OPEN ACCESS

Depression and Suicide Analysis Using Machine Learning and NLP

To cite this article: Pratyaksh Jain *et al* 2022 *J. Phys.: Conf. Ser.* **2161** 012034

View the [article online](#) for updates and enhancements.

You may also like

- [Chernobyl cleanup workers from Estonia: follow-up for cancer incidence and mortality](#)
Kaja Rahu, Anssi Auvinen, Timo Hakulinen et al.
- [Research Progress of University Psychology Based on Big Data-----discovery, challenges and opportunities](#)
Liu Yun
- [Used of Motion Graphics to Create Awareness on Handling Stress](#)
Nur Azila Azahari, Wan NorAshiqin Wan Ali, Tengku Kastriafuddin Shah Tengku Yaakob et al.



The Electrochemical Society
Advancing solid state & electrochemical science & technology

243rd ECS Meeting with SOFC-XVIII

More than 50 symposia are available!

Present your research and accelerate science

Boston, MA • May 28 – June 2, 2023

[Learn more and submit!](#)

Depression and Suicide Analysis Using Machine Learning and NLP

Pratyaksh Jain¹, Karthik Ram Srinivas² and Abhishek Vichare³

^{1,2}Student of Department of Computer Engineering, Mukesh Patel School of Technology Management and Engineering, NMIMS University, Mumbai, India

³Faculty of Department of Computer Engineering, Mukesh Patel School of Technology Management and Engineering, NMIMS University, Mumbai, India.
vichare1@gmail.com

Abstract. Depression is a common type of mental illness that can impair performance and lead to suicide ideation or attempts. Traditional techniques used by mental health experts can assist in determining an individual's type of depression. Machine learning and NLP were used to understand how to predict posts that indicate depression in people and their accuracy. For this work, we have used a dataset from reddit. Reddit is an ideal destination to use as a supplement to the traditional public health system because of its punctuality in exchanging ideas, versatility in presenting emotions, as well as compatibility to use medical terms. We examined the comments and posts about suicidal ideation. We used NLP to gain a better understanding of interdisciplinary fields which are related to suicide. We discovered two help groups for depression and suicidal thoughts: r/depression and r/SuicideWatch. The famous “SuicideWatch” subreddit is commonly used by people who have thoughts of suicide and gives significant signals for suicidal behavior. A brief scan through the articles discloses that the subreddits are legitimate online spots to seek assistance and provide honest text data about people’s mental state. We have used multiple ML algorithms such as Naïve Bayes, SVM. To address the research problem, we have considered two subreddits that provided us with appropriate information to track people at risk. We achieved results of 77.29 % accuracy and 0.77 f1-score of Logistic Regression, 74.35 % accuracy and 0.74 f1-score of Naïve Bayes, 77.120% accuracy and 0.77 f1-score of Support Vector Machine, 77.298% accuracy, and 0.77 f1-score of Random Forest.

1. Introduction

Depression Disorder is a severe and widespread mental illness characterized by an excessive sense of pessimism and despair. Depression, in its most severe manifestations, can have a significant impact on human performance as well as human life. The severity of depression varies from person to person; what most cases have had familiar is a lack of motivation to attempt almost anything, such as tasks they once loved hugely. Suicide ideation and a loss of desire to quantify further are two of the deadliest manifestations of depression. According to psychologists, taking into account depression occurs after two weeks of the existence of its symptoms.

Suicide is defined as a fatal self-harming act committed with the intent to perish. Suicide is the world's 13th major cause of death, accounting for five percent of all demise, according to WHO, with nearly a million people dying as a result worldwide due to suicide each year. Suicide affects all age groups, but worldwide, rates clearly increase with rising age.

There are many reasons people attempt suicide like:



Content from this work may be used under the terms of the [Creative Commons Attribution 3.0 licence](https://creativecommons.org/licenses/by/3.0/). Any further distribution of this work must maintain attribution to the author(s) and the title of the work, journal citation and DOI.

1. **Mental Illness:** Depression can cause people to experience excruciating mental anguish and a sense that all hopes have been dashed, leaving them unwilling to see any way out of the distress except to take their own life.
2. **Traumatic Stress:** An individual that has experienced traumatic events, such as sexual violence, Victims of sexual, physical assault, or mental anguish are more likely to commit suicide several years after the mental anguish has passed because the mental anguish causes feelings of guilt and despair that could result in suicide.
3. **Use of Substances:** Alcohol, drugs also impact a suicidal individual, making them more irrational to act on their impulses than they'd be if they were sober. Individuals with depression and other psychiatric disorders have a risk for substance use and alcohol use dysfunction. When you combine these factors, the dangers multiply.
4. **Fear of Loss:** When faced with a failure or the anxiety of a loss, a person may decide to attempt suicide. These situations can include. These situations can include
 - a) Academic failure
 - b) Bullying, shaming, or humiliation, including cyberbullying
 - c) Financial problems
 - d) Loss of social status

In this paper, we commence with difficulties in processing social media text, which is illustrated with examples of suicidal ideation and similar mental illnesses. As mentioned in the related research, we have discussed prior studies conducted that are relevant to our work. We then provide a thorough analysis of the study which was used in our approach. Ultimately, we outline the method of analysis for meeting the problems outlined in the situation. We conclude by highlighting the broader scope of our research.

2. Literature Review

1. Mining social media data (for example, Twitter, Facebook, and Reddit) for health-related information has justly received a lot of attention recently. There have been streamlined efforts to identify health-related information in social data; also there are results in specific subdomains such as pharmacology, disease surveillance, mental health, and substance abuse supervising. More recent results can be found in a brief survey conducted by Paul et al.
2. S.S Priyanka et al. focused on identifying the primary factors influencing the suicide rate in specific parts of India. Pearson correlation and OLS regression resulted in 0.998 R-squared value and 0.991 Adjusted R-Squared value.
3. It's been identified that social networking sites can perform both positive and negative functions in suicide-related studies associated with social media. Won et al. conducted an intriguing survey that demonstrated that social media parameters based on blog posts are highly associated with national suicide rates in Korea. Depending on the Twitter data set, Jashinsky et al. later encountered a similar result on the US population. Burnap et al. recently developed text classifiers that recognize tweets about suicidal behavior. They also investigated the follower-friend networks of tweeters who posted tweets expressing suicidal behavior.
4. Alambo et al., proposed their approach forward towards a possible answer for a computerized suicide risk-elicitation framework. Their two-pronged method takes advantage of 1) clustering and 2) sequence modeling techniques. Furthermore, they compile a database of suicide threads along with their level of risk.
5. J. H. K. Seah et al. perform social sensing on Reddit-sourced digital traces; they examine the posts and comments about depression and suicide. They use NLP to gain a better understanding of diverse aspects that are related to suicide. Their research shows that using an approach based on data mining enables accurate as well as the computerized method to detect suicide on media platforms.
6. G. M. Lin et al. employed Based on six critical psychological stress domains; ML algorithms were used to recognize the existence of suicidal self-harm thoughts in military males and females. All six machine learning methods are more than 98 percent accurate. The

- multilayer perceptron and support vector machine, among others, get the most precise forecasts of suicidal ideation nearly 100 percent of the time.
7. S. Fodeh et al. proposed a suicide clinical trials may benefit from ML models. They downloaded 12,066 series of posts from 3,873 users using keywords from Jashinsky et al. ideation's tracking framework via Twitter's Android app. Participants were allocated "HighRisk" or "at risk" keywords relying on their use of suicidal ideation aspects, and optimization classes have been used to recognize implied suicidal behavior potential risks between many users, that were then used to identify and classify users as "HighRisk" or "at risk." Amongst these algorithms used were Semantic Analysis, Latent Dirichlet allocation, and Non-linear Programming.
 8. F. Chiroma et al. investigates the computational efficiency classifiers in detecting suicidal social media posts. The test employed four well-known machine classification techniques machine; the experiment produced an F-measure ranging from 0.346 to 0.778 for suicidal interaction, with the Decision Tree algorithm achieving the highest efficiency.
 9. X. Huang et al. present a new dataset containing 130 Chinese social media profiles of suicide victims. They discovered spectral patterns that correspond to how people perceive themselves in the weeks and months before suicide: an excessive number of posts, a rise in suicidal thoughts, and a rise in bearish attitudes in the penultimate weeks and days before suicide.
 10. De Choudhury et al. addressed the sequential problem of recognizing changes to suicidal behavior from Reddit mental health discussion. They designed the move to ideation in terms of users posting in other mental health subreddits for an amount of time, accompanied by a post in the SW subreddit for a consequent amount of time. They have a 77.5 percent success rate in categorizing users who have developed suicidal thoughts.

Table 1. Literature review

Publications	Inputs	Methods
S. A. S. A. Kulasinghe et al. [2]	Text	Conversational Model, Voice Analysis, Feature Extraction, Neural Network
S.S Priyanka et al. [3]	Text	Pearson correlation, OLS regression
A. Alambo et al. [4]	Novel question-answering	Semantic clustering, Sequence to sequence modeling techniques
J. H. K. Seah et al. [5]	Reddit sourced Digital Traces	Maximum Entropy, Conditional Random Field, Hybrid Machine Learning
G. -M. Lin et al. [7]	Text	multilayer perceptron
S. Fodeh et al. [8]	Ideation tracking structure to download tweets via Twitter's API	Latent Semantic Analysis, Latent Dirichlet allocation, Non-linear Programming
X. Huang et al. [10]	Twitter Depressed Users data	Longitudinal text analysis
J. Li et al. [11]	Texts from Blogs	emotion classification, emotion recognition, Affective computing

A. Carrillo-Morales et al. [12]	Text	Text Processing, text analysis
Y. Tai et al. [13]	Questionnaire	Risk analysis, Neural Networks , ROC , Radial basis function training, Past histories prediction
J. Baek et al. [14]	Text	Context Deep Neural Network Model , Multiple regression , Context-DNN model , Regression analysis
M. M. Tadesse et al. [15]	Reddit posts	Multilayer Perceptron classifier, SVM , Naïve Bayes
A. M. Schoene et al. [16]	Text	Recurrent neural networks, Task analysis , Text Classification
A. Seal et al. [17]	Electroencephalogram (EEG) data	Convolutional neural nets , Support Vector Machine , Sensitivity analysis
S. B. Hassan et al. [18]	Text	Pattern classification, Support vector machines, Text analysis
M. Trozsek et al. [19]	User-level linguistic metadata	Convolutional neural network , ERDE score , Text sequences , Word Embeddings
S. Tokuno et al. [20]	Voice data from depressed military	Speech recognition , Emotion Recognition
V. R. Chiranjeevi et al. [21]	Text	Feature extraction , Self-Organizing Map (SOM) , Neural Networks

In comparison to all these endeavors, we concentrate on the combined task of recognizing useful comments to SW subreddit posts. The objective is to predict helpful comments as well as analyze and interpret results in order to gain insights into communication techniques for online responses to suicidal social media posts from users.

3. Methodology

We have followed the methodology of collecting data, then pre-processed it, then trained the model and finally validated the model.



Figure 1. Figure indicating the methodology of data collection and training.

Data collection:

The dataset consists of over 60000 individual data points. Classified as depression vs. suicide, data was collected from two subreddits. Some of the characteristics and properties which make the dataset ideal for use in our model are the active postings in them, low troll rate and low spam rates and more textual tweets with very few images posted. Less memes also help in cleaning the data. With our research, we are trying to identify the distinction between language used by a person experiencing clinical depression thoughts and language used by an individual at risk of suicide, which will be beneficial for counsellors and mental health professionals.

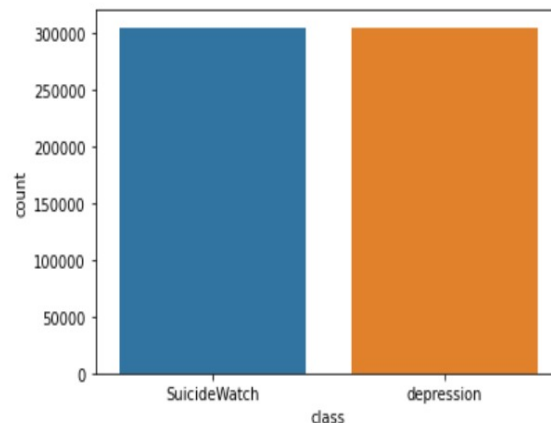


Figure 2. Graph comparing count and class of the model

Data pre-processing

We converted the text in lower casing, removed the punctuation, removed numbers, tokenized the data, removed the stop words, stemmed the data and then finally passed it into machine learning model.

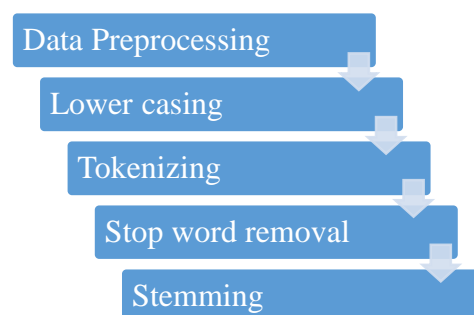


Figure 3. Figure indicating steps in data pre-processing.

Lower casing:

Converting a word to lower case, words like 'Depressed' and 'depressed' mean the same, but when not converted to the lower case, those two are represented as two different words in the vectors.

Tokenization:

Tokenization is a process of converting text to a smaller unit called tokens; tokens can be of 3 types: words, characters, and ngrams. We have broken down our text data into words.

Stop words:

Stop words are words in the language that do not add significant meaning to a sentence. Common stop words are he, his, your, does, do, if etc. they can be easily ignored without losing the meaning of the sentence. By removing such words, we reduce the overall data size without reducing any valuable information.

Stemming:

Stemming is the process of removing a part of the word, or reducing it to its stem or root.

ML Algorithms Applied

Logistic regression:

It predicts a binary outcome based on a set of independent variables. In contrast to linear regression, which yields constant number values, logistic regression yields a probability value that can be linked to two distinct classes using the logistic function.

Naïve bayes:

It uses bayes theorem to classify objects. Whenever we are dealing with a large amount of data, naïve bayes is the preferred solution. This provides very good results when it comes to NLP activities, such as sentimental analysis as it's a fast and uncomplicated classification algorithm.

Support vector machine:

SVM generates a linear algorithm that creates a function to maximize the distance between classes. Class function is created using instance data at the edge of the class. The data points lying closest to the decision boundary are the most difficult to identify, and they have a direct impact on the optimum position of the decision surface. SVM algorithm can achieve this by creating a function that maximizes the distance.

Random forest:

It is made up of many decision-making trees; each provides a class forecasting. The model prediction is chosen from the class with the most votes. The concept is simple but robust, and it represents the wisdom of the masses.

Random forest also de-correlates the independently constructed trees by selecting a subset of the available features (predictors) at random to build the tree. This procedure eliminates the possibility of the independently constructed trees being highly correlated due to one or two extremely powerful predictors.

4. Results and findings

Logistic Regression:

Table 2. Table indicating the precision, recall, and f1-score for logistic regression.

	Precision	Recall	f1-score
1	0.80	0.79	0.79
2	0.79	0.80	0.79

The results show precision 0.80, recall of 0.79, f1-score 0.79 and accuracy of 0. 79. Logistic regression has had the highest accuracy along with random forest, followed by Support Vector Machine, Naïve Bayes.

Naïve Bayes:

Table 3. Table indicating the precision, recall, f1-score for Naïve Bayes.

	Precision	Recall	f1-score
1	0.78	0.72	0.75
2	0.74	0.80	0.77

The results show precision 0.78, recall of 0.72, f1-score 0.75, and accuracy of 0.7586808212887849. Naïve Bayes has had the lowest accuracy among the algorithms followed by SVM then Logistic Regression and Random Forest.

Support Vector Machine:

Table 4. Table indicating the precision, recall, and f1-score for Support Vector Machine.

	Precision	Recall	f1-score
1	0.77	0.78	0.77
2	0.77	0.76	0.77

The results show precision 0.77, recall of 0.77, f1-score 0.77, support of 3938 and accuracy of 0.77120. Support Vector Machine has had the second-highest accuracy higher than Naïve Bayes but less than Logistic Regression and Random Forest.

Random Forest:

Table 5. Table indicating the precision, recall, and f1-score for Random Forest.

	Precision	Recall	f1-score
1	0.78	0.77	0.77
2	0.77	0.77	0.77

5. Conclusion

Depression and suicide analysis is considered as a challenging and complex task. In this paper, we attempted to detect the existence of depression on the Reddit platform and looked for ways to improve affective efficiency in order to detect depression. We discovered a stronger link between depression and language using text classification techniques and NLP. We looked at how single feature and cumulative feature sets performed in measuring depression symptoms using various text classifying methods.

The “SuicideWatch” subreddit is used by people who have thoughts of suicide and gives great signals for suicidal behavior. Reddit is an ideal destination to use as a supplement to the traditional public health system because of its punctuality in exchanging ideas, versatility in presenting emotions, as well as compatibility to use medical terms. The model's effectiveness can be seen with 77.29 % accuracy and 0.77 f1-score of Logistic Regression, 74.35 % accuracy and 0.74 f1-score of Naïve Bayes, 77.12% accuracy and 0.77 f1-score of Support Vector Machine, 77.298% accuracy, and 0.77 f1-score of Random Forest.

Even though techniques used in this work perform sufficiently well, but even more research require in this area. We presume that this study will help to lay the groundwork for new mechanisms that will be used to approximate depression and related factors in various fields of health. People suffering from mental illnesses may benefit from being more proactive in their recovery.

Identifying the instant suicide risk is a challenging but possibly lifesaving task. The existing accepted indicators do not satisfactorily tackle the topic of Suicide risk. Although the presence of suicidal ideation indicates a high risk of suicide, it is clear that many patients deny having one. Suicidal thoughts, intentions, and information from significant others should be sought and considered legitimate even after the patient denies it.

6. Future Work

Our structure can be widened to solve more significant healthcare problems encompassing multimodal data. It can even be used in conjunction with smart virtual assistants to lower the risk of self-harm in patients. Some points to improve in our present system

1. Collaborating with local groups to integrate digital footprints with offline data (e.g., suicide hotlines, counseling)
2. Extending our research beyond Reddit and into other common public forums.

Further DL algorithms such as RNN (Recurrent neural networks), LSTM (Long Short-Term Memory), ANN (Artificial Neural Network) can be implemented to improve the research. A further direction for future research is to take into account other data characteristics such as comments, memes, reposts, and likes. It would also be interesting to investigate whether such conversations affect the direction of a user's suicide sequence of events.

7. References

- [1] Benjamin L. Cook, Ana M. Progovac, Pei Chen, Brian Mullin, Sherry Hou and Enrique Baca-Garcia 2016. Novel Use of Natural Language Processing (NLP) to Predict Suicidal Ideation and Psychiatric Symptoms in a Text-Based Mental Health Intervention in Madrid. *Computational and Mathematical Methods in Medicine*, vol. 2016, Article ID 8708434, 8 pages. <https://doi.org/10.1155/2016/8708434>
- [2] S. A. S. A. Kulasinghe, A. Jayasinghe, R. M. A. Rathnayaka, P. B. M. M. D. Karunaratne, P. D. Suranjini Silva and J. A. D. C. Anuradha Jayakodi 2019. AI-Based Depression and Suicide Prevention System. *International Conference on Advancements in Computing (ICAC)*, pp. 73-78, DOI: 10.1109/ICAC49085.2019.9103411.
- [3] S. S. Priyanka, S. Galgali, S. S. Priya, B. R. Shashank and K. G. Srinivasa 2016. Analysis of suicide victim data for the prediction of number of suicides in India. *2016 International Conference on Circuits, Controls, Communications and Computing (I4C)*, pp. 1-5, DOI: 10.1109/CIMCA.2016.8053293.
- [4] Alambo, A., Gaur, M., Lokala, U., Kursuncu, U., Thirunarayan, K., Gyrard, A., Sheth, A., Welton, R. S., Pathak, J. 2019. Question Answering for Suicide Risk Assessment Using Reddit. *2019 IEEE 13th International Conference on Semantic Computing (ICSC)*, pp. 468-473, DOI: 10.1109/ICOSC.2019.8665525.
- [5] J. H. K. Seah and K. Jin Shim 2018. Data Mining Approach to the Detection of Suicide in Social Media: A Case Study of Singapore. *2018 IEEE International Conference on Big Data (Big Data)*, pp. 5442-5444, DOI: 10.1109/BigData.2018.8622528.
- [6] A. Carrillo-Morales and A. Curiel 2019. Advances Towards the Identification of Mental Disorders Associated with Suicide through Text Processing. *2019 International Conference on Inclusive Technologies and Education (CONTINUE)*, pp. 121-1217, DOI: 10.1109/CONTIE49246.2019.00031.

- [7] G. -M. Lin, M. Nagamine, S. -N. Yang, Y. -M. Tai, C. Lin and H. Sato 2020. Machine Learning-Based Suicide Ideation Prediction for Military Personnel. *IEEE Journal of Biomedical and Health Informatics*, vol. 24, no. 7, pp. 1907-1916, July 2020, DOI: 10.1109/JBHI.2020.2988393.
- [8] S. Fodeh et al. 2019. Using Machine Learning Algorithms to Detect Suicide Risk Factors on Twitter. *2019 International Conference on Data Mining Workshops (ICDMW)*, pp. 941-948, DOI: 10.1109/ICDMW.2019.00137.
- [9] F. CHIROMA, H. LIU and M. COCEA 2018. Text Classification For Suicide-Related Tweets. *2018 International Conference on Machine Learning and Cybernetics (ICMLC)*, 2018, pp. 587-592, DOI: 10.1109/ICMLC.2018.8527039.
- [10] X. Huang, L. Xing, J. R. Brubaker and M. J. Paul 2017. Exploring Timelines of Confirmed Suicide Incidents Through Social Media. *2017 IEEE International Conference on Healthcare Informatics (ICHI)*, pp. 470-477, DOI: 10.1109/ICHI.2017.47.
- [11] J. Li and F. Ren 2008. Emotion recognition from blog articles. *2008 International Conference on Natural Language Processing and Knowledge Engineering*. pp. 1-8, DOI: 10.1109/NLPKE.2008.4906757.
- [12] A. Carrillo-Morales and A. Curiel 2019. Advances Towards the Identification of Mental Disorders Associated with Suicide through Text Processing. *2019 International Conference on Inclusive Technologies and Education (CONTINUE)*, pp. 121-1217, DOI: 10.1109/CONTIE49246.2019.00031.
- [13] Y. Tai and H. Chiu 2007. Artificial Neural Network Analysis on Suicide and Self-Harm History of Taiwanese Soldiers. *Second International Conference on Innovative Computing, Information and Control (ICICIC 2007)*, pp. 363-363, DOI: 10.1109/ICICIC.2007.186
- [14] J. Baek and K. Chung 2020. Context Deep Neural Network Model for Predicting Depression Risk Using Multiple Regression. *IEEE Access*, vol. 8, pp. 18171-18181, DOI: 10.1109/ACCESS.2020.2968393.
- [15] M. M. Tadesse, H. Lin, B. Xu and L. Yang 2019. Detection of Depression-Related Posts in Reddit Social Media Forum. *IEEE Access*, vol. 7, pp. 44883-44893, DOI: 10.1109/ACCESS.2019.2909180.
- [16] A. M. Schoene, A. Turner, G. R. De Mel and N. Dethlefs. Hierarchical Multiscale Recurrent Neural Networks for Detecting Suicide Notes. *IEEE Transactions on Affective Computing*, DOI: 10.1109/TAFFC.2021.3057105.
- [17] A. Seal, R. Bajpai, J. Agnihotri, A. Yazidi, E. Herrera-Viedma and O. Krejcar 2021. DeprNet: A Deep Convolution Neural Network Framework for Detecting Depression Using EEG. *IEEE Transactions on Instrumentation and Measurement*, vol. 70, pp. 1-13, Art no. 2505413, DOI: 10.1109/TIM.2021.3053999.
- [18] S. B. Hassan, S. B. Hassan and U. Zakia 2020. Recognizing Suicidal Intent in Depressed Population using NLP: A Pilot Study. *11th IEEE Annual Information Technology, Electronics and Mobile Communication Conference (IEMCON)*, pp. 0121-0128, DOI: 10.1109/IEMCON51383.2020.9284832.
- [19] M. Trotzek, S. Koitka and C. M. Friedrich 2020. Utilizing Neural Networks and Linguistic Metadata for Early Detection of Depression Indications in Text Sequences. *IEEE Transactions on Knowledge and Data Engineering*, vol. 32, no. 3, pp. 588-601, 1 March, DOI: 10.1109/TKDE.2018.2885515.
- [20] S. Tokuno et al. 2011. Usage of emotion recognition in military health care. *Defense Science Research Conference and Expo (DSR)*, pp. 1-5, DOI: 10.1109/DSR.2011.6026823.
- [21] V. R. Chiranjeevi and D. Elangovan. Surveillance Based Suicide Detection System Using Deep Learning. *2019 International Conference on Vision Towards Emerging Trends in Communication and Networking (ViTECoN)*, pp. 1-7, DOI: 10.1109/ViTECoN.2019.8899360.

- [22] A. Zahura and K. A. Mamun 2020. Intelligent System for Predicting Suicidal Behaviour from Social Media and Health Data, *2nd International Conference on Advanced Information and Communication Technology (ICAICT)*, pp. 319-324, DOI: 10.1109/ICAICT51780.2020.9333463.
- [23] X. Meng and J. Zhang 2020. Anxiety Recognition of College Students Using a Takagi-Sugeno-Kang Fuzzy System Modeling Method and Deep Features. *IEEE Access*, vol. 8, pp. 159897-159905, 2020, DOI: 10.1109/ACCESS.2020.3021092.

Acknowledgment

We would like to thank our university NMIMS (Mumbai Campus), our college MPSTME and Computer Department for giving us the opportunity of writing this paper.