# Social Media Interaction-Based Mental Health Analysis with a Chat-Bot User Interface

**Aliyah Kabeer, Paul John, Serena A. Gomez, Pooja Agarwal, and U. Ananthanagu**

**Abstract**  In recent times, social media has played a major role in shaping the state of mind of young adults. This means that the content that is interacted with online may subconsciously have a significant effect on one's mental health. This paper presents a novel approach to detect the mental health state of a user by analyzing their online activity over a period of time and generating a report indicating the same. The user is then allowed to ask any clarifying questions on the generated report, along with general queries on mental health via a chat-bot interface. In order to classify the mental state of the user based on the kind of content posted or interacted with by them, their tweets are scraped and feature vectors of the same are generated. Supervised machine learning algorithms like Support Vector Machines (SVM) and neural network-based models like Long Short-Term Memory (LSTM) are compared for their performance on prediction. A transfer learning approach is also attempted and gives promising results in predicting the classes of the tweets. Natural Language Processing techniques such as question similarity and extractive summarization are utilized in building the chat-bot framework.

**Keywords**  Social media · Mental health · Report · Chatbot · Question similarity · Twitter · Tweets

A. Kabeer (✉) · P. John · S. A. Gomez · P. Agarwal · U. Ananthanagu
Department of Computer Science and Engineering, PES University, Bangalore, India
e-mail: aliyahk8888@gmail.com

S. A. Gomez
e-mail: serenagomez@pesu.pes.edu

P. Agarwal
e-mail: poojaagarwal@pes.edu

U. Ananthanagu
e-mail: ananthanagu@pes.edu

# 1   Introduction

Mental health is a state of emotional and cognitive well-being which affects a person's thought process, feelings, and actions while coping with the normal stresses of everyday life. In today's world, keeping track of ones own mental health is of utmost importance, given the hectic and fast-paced lives that we live. Social media platforms such as Twitter have become a popular platform on which people of all age groups express their emotions, interests, likes, dislikes, and day-to-day activities. Since an average of 500 million tweets are tweeted every single day, it seems necessary to keep track of a person's mental health to ensure that his or her tweets do not indicate a declining mental health and to alert that person if necessary. An intelligent chat-bot is a computer program that can maintain a conversation or answer questions based on the context of the input provided. Considering that most states of mental health can best be dealt with by having a conversation, a chat-bot is the ideal program to integrate into the application.

Through the process of multi-class sentiment analysis, the proposed framework can classify the tweets of a person into various categories such as normal, anxious, stressed, suicidal, and lonely. Sentiment analysis is performed using various machine and deep learning algorithms such as Recurrent Neural Networks, Support Vector Machines, and Logistic Regression. By fixing a certain threshold value, each tweet is classified into one of these classes and the model generates a report based on these tweets over fixed periods of time.

Since the mental health of a person is an extremely delicate and sensitive topic, it would only be normal for the user to be concerned about the generated report. The proposed framework includes an intelligent chat-bot that is trained to answer questions related to the generated report and mental health in general, along with the ability to maintain a conversation with the user, thus helping the user feel safe, aware, and taken care of. It is a program that matches the input with the most appropriate answer or response. It will be made aware to the user that the chat-bot in no way will be able to provide a diagnosis for the predicted report, but rather make suggestions and indicate which tweets were classified what. The chat-bot further provides a more detailed answer for any mental health query by finding the three most similar questions and summarizing their corresponding answers. By performing the summarization of answers, it is ensured that the chat-bot responds in a more human-like manner, by generating a different answer if the question is framed differently, instead of retrieving static responses from the dataset every time.

# 2   Scope and Motivation

Social media being easily accessible and convenient to use allows individuals to keep in touch without actually having to be physically present. However, this incessant hyper-connectivity can negatively impact one's mental state and trigger impulse

control problems. Studies have shown the existence of a strong link between the usage of social media and an escalated risk for anxiety, stress, loneliness, and suicidal ideations. Moreover, the subconscious effects of social media on mental health might not be obvious immediately but might lead to a snowball effect over time that might require constant and immediate attention. The only way to keep a user's online activity in check is by monitoring it, which include but are not limited to the content that they post, like, re-share, save, comment on, and follow. Virtual assistants are devices that have access to a user's online activity. Having a plugin that tracks this usage over periods of time and classifies which content might be negatively impacting their moods or mental state can allow them to be more self-aware and steer clear of such negative content. Thus, one application of the proposed framework involves posing as a virtual assistant plugin that will allow individuals to keep track of their social media usage and assess their mental state along the way.

## 3   Related Work

Of late, mental health analysis and prevention of the onset of mental health illnesses have been a trending research topic and are gaining popularity among researchers due to its wide scope. Moreover, sentiment analysis techniques and the use of Natural Language Processing to extract semantic information from text have been recently extended to the field of mental health since the rise in social media usage over the previous few years has been prominently high. Yatapala et al. [1] present a machine learning approach using artificial neural networks to detect suicide ideation or thoughts and identify patterns in suicidal text by analyzing tweets of the user. The paper also compares the performance of the Word2Vec and TF-IDF vectorizers to generate the feature vectors for their model. They concluded by proving that TF-IDF has worked better and obtained a testing accuracy of 89% . Tiwari and Verma [2] talk about the rapid growth of sentiment analysis on text data provided by users and developed a model to classify new tweets into fine grained emotions such as angry, love, happy, boredom, sad, fun, surprise, neutral, empty, relief, enthusiastic, worry, and hate. They worked with 40,000 tweets, using sentence-level analysis, split into an 80:20 ratio for training and testing, and used three machine learning prediction algorithms, Decision Tree, Support Vector Machines, and Random Forests, which all recorded an accuracy of greater than 91%. Jain et al. [3] also talk about suicidal ideation and depression detection using Twitter data and questionnaires. They utilized classification algorithms like Random Forest, XGBoost, Support Vector Machine, and Logistic Regression to classify the severity of the detected depression into five stages. Although they analyze the user's Twitter activity, they do not provide the user any means to identify the content they interacted with that might have indicated their mental state at the time.

Khariya and Khodke in [4] utilized the Twitter API to fetch tweets and classify them into positive, negative, or neutral sentiments by using algorithms such as K-Nearest Neighbors and Naive Bayes. Although high accuracies were obtained, these

models are not context-aware and do not consider the semantic information of entire sequences of text. Dinakar et al. [5] talk about a way of performing sentiment analysis on Twitter data to extract contextual polarity that is to determine if the given content is predominantly negative or not. This is done by the preparation of a lexicon wherein each word is given a value that indicates its contextual polarity. On discovering negative polarity, the most common tags or terms used by the user are further isolated by performing clustering of data. The result is provided in the form of the percentage of tweets that are either positive or negative and also the visualization of the same on a histogram. One of the key limitations found was the existence of only a broad classification of the tweets into positive or negative classes.

Conversational and dialogue agents like chatbots are often associated with mental health-related applications owing to the fact that bots can provide a more reliable and unbiased channel for people to talk to and express their feelings. The earliest mental health chat-bot therapist ELIZA utilized pattern matching and substituted methodology to make users believe that they were talking to a human and give them a safe space to share their feelings. Amer et al. in [6] built a domain-specific chat-bot framework to perform the task of question answering by finding the answer to a user's query by taking a reference context passage as input and locating its answer. This limited its application and required a more enriched dataset to improve the model accuracy and robustness. Hwerbi and Khouloud [7] built a chat-bot containing various components—the ontology, web scraping module, database, state machine, keyword extractor, a trained chat-bot, and finally, the user interface. Here, they develop an ontology by acquiring knowledge via web scraping, defining competency questions, concepts necessary to answer predefined questions, and properties of concepts, individuals, and between individuals. This was followed by the implementation of a state machine, after which they make use of a keyword extractor—Rapid Automatic Keyword Extraction (RAKE) that uses a comprehensive list of phrase delimiters and stop words in order to identify the most relevant words or phrases in a given text. They then integrate a trained chatbot—ChatterBot which is a conversational dialog engine that was built in Python language and has the ability to generate meaningful responses based on collections of known conversations. Finally, they visualize the data on a map and add a user interface. Lalwani and Rathod [8] implemented an AIML-based chatbot that stores the question and answer pairs in AIML files, where the input is matched with the questions using semantic similarity algorithms, and the appropriate answer is returned as a response. They performed lemmatization and POS tagging using WordNet, which is a lexical database for English. Their paper focused on implementing a chatbot interface for a College Inquiry System to answer college-related questions, retrieve general information of the students, and get information on upcoming events.

# 4 Proposed Methodology

The proposed methodology includes generating a mental health report based on the users' social media activity and providing them with the ability to ask any questions regarding the generated report and regarding general mental health-related questions via a chat user interface. Figure 1 shows the proposed methodology of the application. The social media corpus comprises a cleaned dataset of tweets categorized into five labels which serve as the training data for the proposed model. The tweets are then preprocessed and their features are extracted in order to generate the feature vectors. Finally, the preprocessed data is fed into a machine or deep learning model for training purposes. This trained model is used to predict target labels for real-time tweets interacted with or posted by the user and present the results in the form of a pie chart. A mental health FAQ and report-related corpus are then used to train the chat-bot component of the user interface, which can answer the user's queries post-report generation. Figure 2 shows the logic flow followed by the bot for answer generation. The detailed design of each module is as given below:

## 4.1 Social Media Activity Corpus

The social media corpus comprises a cleaned dataset of tweets categorized into five labels, namely—anxious, lonely, stressed, normal, and suicidal, each indicating the probable mental state of the user at the time of tweeting or interacting with tweets. This data serves as the training data for the proposed mental health assessment model.

## 4.2 Mental Health Assessment Model

**Data Collection** Twitter is being one of the most popular and widely used social media platforms globally, and the number of tweets generated every second is unbelievably high. This serves as a continuous source of information for analyzing a user's
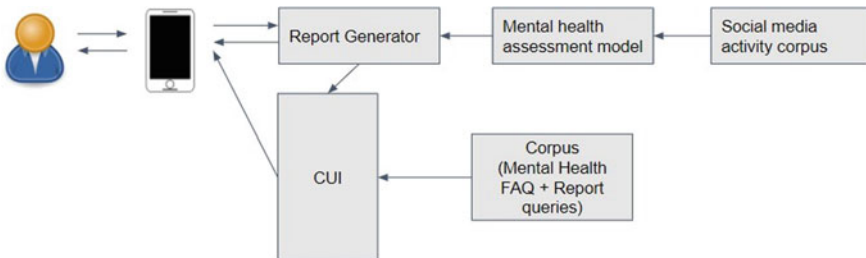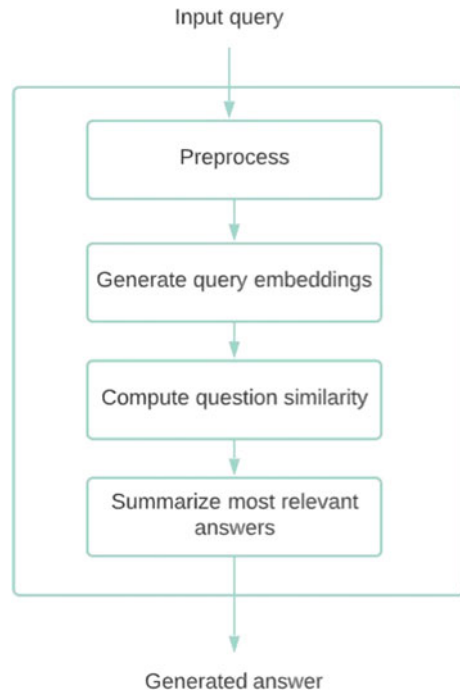


**Fig. 1** High-level architecture

**Fig. 2** Logic flow for
answer generation

Input query

Preprocess

Generate query embeddings

Compute question similarity

Summarize most relevant
answers

Generated answer

state of mind over a period of time. The social media activity corpus contains labeled
tweets which will be used as the training dataset.

**Data Preprocessing** Data collected from the real world may be inconsistent, noisy,
and incomplete. For this reason, data preprocessing is a crucial part of nearly every
machine learning process to ensure that the model is trained on consistent, reliable
data. It further helps in analyzing and visualizing the distribution of data to understand
it better. First, the necessary libraries for preprocessing are imported, following which
the tweets are tokenized. Next, the stop words, numbers, punctuation, and special
characters are removed. This forms a key part of preprocessing as tweets are generally
comprised various special characters like hashtags, hyperlinks, urls, and numbers,
which might not be relevant to the application and might have an adverse negative
effect on the training process. The tokens are then lemmatized in order to get to their
base form.

**Feature Extraction** The feature vectors are formed using TF-IDF or Term
Frequency-Inverse Document Frequency, which is a technique to extract weighted
word count features. The cleaned data is first split into a 80–20 train–test split ratio,
wherein 80% of the data is reserved for training purposes, and the remaining portion
for testing. Random shuffling of the partitioned data ensures equal representations
of classes. This data is then fed into the TF-IDF model which produces a vectorized
form of the tweets.

**Machine Learning Models** In order to classify the mental state of a user based on their tweets, sentiment analysis and text classification algorithms are utilized. Since the training dataset is labeled, supervised machine learning algorithms such as Support Vector Machines (SVM) and neural network-based algorithms like Long Short-Term Memory (LSTM) are trained and their performances are compared. The details of the two models are as given below:

- *Support Vector Machines:* Support Vector Machine (SVM) is a type of supervised machine learning algorithm that maps data to higher dimensions in order to perform linear or nonlinear classifications with the help of the kernel trick. The linear SVM classifier uses a linear kernel function and is particularly used for text classification since it was found that text is most often linearly separable. Moreover, since text has a lot of features, mapping it to a higher dimension will not help solve the problem. The SVM classifier obtained an accuracy of 87.87% on the testing data. The classification report of the aforementioned model is given in Fig. 3, where it is seen that the model has a high precision, recall, F1-score, and support for all classes except for those classes labeled "lonely." This anomaly could be because there is relatively lesser amount of data available for that particular class in the training dataset.
- *Long Short-Term Memory (LSTM)* LSTM is a kind of Recurrent Neural Network that can grasp sequences and long-term dependence. They are largely used for different tasks like sentiment analysis, text classification, speech recognition, etc. LSTMs have the advantage of retaining only useful information and discarding the remaining part, and this property allows it to model the context of various words in the text. Traditional machine learning models follow a statistical approach to find the probability of a sentence belonging to a specific class by training on each word vector individually. This would require custom features or a bag-of-words representation to be created. On the other hand, LSTMs can see the text as an entire sequence and work from there. The LSTM model for the proposed architecture

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.89 | 0.89 | 0.89 | 1844 |
| 1 | 0.91 | 0.89 | 0.90 | 789 |
| 2 | 0.33 | 0.11 | 0.17 | 9 |
| 3 | 0.84 | 0.87 | 0.85 | 1146 |
| 4 | 0.87 | 0.75 | 0.81 | 81 |
| accuracy |  |  | 0.88 | 3869 |
| macro avg | 0.77 | 0.70 | 0.72 | 3869 |
| weighted avg | 0.88 | 0.88 | 0.88 | 3869 |

```
Accuracy:
0.8787800465236495
```

**Fig. 3** Classification report for Support Vector Machine

uses only one LSTM hidden layer followed by a dense layer. The outputs are then passed into a ReLU activation function, following which a regularization technique, Dropout, is applied. The final output layer is followed by a Softmax activation function that gives the probability of the tweet belonging to each of the five classes. It is compiled using the RMSProp optimizer and the sparse categorical cross-entropy loss function. The model gives an accuracy of 89% on the testing data, on training for 6 epochs. Figure 4 shows a summary of the LSTM model and its various layers. On comparing the two models on the test data, it is observed that the accuracy of the LSTM is slightly better than that of the SVM model. However, the Support Vector Machines gives a slightly better performance in terms of realistically classifying the real-time tweets.

**Transfer Learning** Transfer learning is a machine learning technique where the knowledge gained when a model is developed for one task is reused as the starting point for a model on a second different task. The proposed framework was also implemented using transfer learning techniques, in order to evaluate and compare the results on multi-class classification of the tweets. The results obtained are compared against those of the machine learning models as seen in Table I. The proposed method uses a sentence transformer library, which makes use of the pre-trained BERT model to compute dense vector representations of the user's tweets. It then makes use of the Optuna model [9], which helps in identifying the best set of hyperparameters that

```
Model: "model"

Layer (type)                 Output Shape              Param #
=================================================================
inputs (InputLayer)          [(None, 500)]             0

embedding (Embedding)        (None, 500, 50)           100000

lstm (LSTM)                  (None, 64)                29440

FC1 (Dense)                  (None, 256)               16640

activation (Activation)      (None, 256)               0

dropout (Dropout)            (None, 256)               0

out_layer (Dense)            (None, 5)                 1285

activation_1 (Activation)    (None, 5)                 0

=================================================================
Total params: 147,365
Trainable params: 147,365
Non-trainable params: 0
_____
None
```

**Fig. 4** Summary of the LSTM model

can be used for the multi-class text classification model to classify the user's tweets. The transfer learning attempt gave a mean accuracy of 80% and shows promising results for future work.

## 4.3 Report Generator

Tweepy is an open-source and easy-to-use Python package that gives access to the Twitter API. This package is used in order to scrape tweets in real-time, given the user's Twitter ID and the start and end dates of the analysis, along with the Twitter developer credentials. The tweets are then cleaned, preprocessed, vectorized, and stored in a data frame for further classification. At the time of report generation, the stored tweets are passed into the mental health assessment model and classified as being normal, stressed, lonely, anxious, or suicidal, each representing the mental state of the user as indicated by the contents of the tweet interacted with or posted by them. The report is presented in the form of a pie chart as can be seen in Fig. 5. A visual representation of the report rather than a textual one ensures that it is easily understandable and helps the user grasp a quick summary of their social media activity at a glance. The colors assigned to the pie chart are indicative of the state of mind that they represent—for instance, yellow being the accepted standard color for happiness is assigned to normal tweets, gray for anxious tweets, and so on. It is to be noted that the results generated are purely a warning for the user that he/she/ they might be feeling that condition and in no way does it diagnose the user as that specific class.

## 4.4 Corpus (Mental Health FAQ + Report Queries)

The mental health FAQ and report query corpora comprise cleaned datasets of mental health-related questions and questions related to the generated report. This data serves as the training data for the conversational user interface (CUI).

## 4.5 Conversational User Interface (CUI)

Once the report is generated, the user is allowed to ask any queries that they may have regarding the results in the report and also queries regarding mental health in general. This task is performed by the chat-bot by first preprocessing the query and then embedding it. This is a key step because if the question is framed differently, the vector representation is still supposed to capture the semantic information in it. For instance, if the user queries "what are the symptoms of mental illness?" or "what are the signs of mental illness?" The bot should be able to vectorize the query in a way
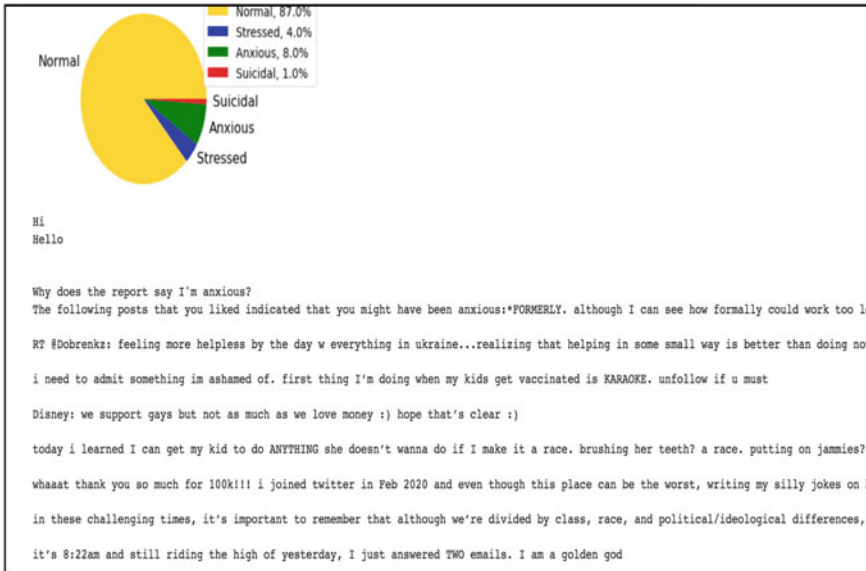
**Fig. 5** Report generation followed by the CUI

that highlights that the key intent of both the queries is the same. This is performed by using the Universal Sentence Encoder. This encoder model was presented in [10] targeting transfer learning to various NLP tasks by generating embeddings in the form of vectors to encode sentences. This performs sentence embedding, which maps the semantic information of the sentences into vectors of real numbers. It uses attention and computes context-aware representations that consider all the other words in the sentence as well—by taking into account both their ordering as well as their identity. Finally, the obtained word representations that are context-aware are converted into a sentence encoding vector of a fixed length by calculating an element-wise sum of the word representations at each position. This helps in understanding the context of the sentences in a more comprehensive way to enhance the process of finding question similarities and retrieving information. The appropriate answer is then retrieved by the bot in the following manner:

- Using the question similarity function to check whether the query asked is a report-related query or not by setting a threshold and checking whether the similarity score of the most similar query is greater than the threshold. If this condition is satisfied, then the answer corresponding to the most similar query is returned.
- Using the question similarity function to check whether the query asked is a mental health-related query or not by setting a threshold and checking whether the average of the similarity scores of the top three most similar queries is greater than the threshold. If this condition is satisfied, then a summarized answer of the top three most similar queries' corresponding to answers is returned. An extractive

summarization technique is used wherein the most important sentences are identified by assigning a score to each sentence based on the importance or frequency of the words present in them. The sentences are then ranked by importance, and the top n sentences with the highest scores are selected to summarize the answer.

- If both the above conditions are not satisfied, then the bot utilizes Python's pretrained chatterbot library for its answer. This bot is trained on numerous English corpora to be able to handle general conversation such as greetings, conversation, emotions, trivia, health, history, humor, politics. This was incorporated in order to ensure that the bot could respond to statements outside the domain-specific queries from the user and carry forward general conversation.

## 5  Results

The mental health status report generated based on the user's social media activity has an accuracy of 87.87% using the Support Vector Machines as the classifier. The chat-bot interface performed relatively well on domain-specific questions like mental health FAQ's and report-related queries. Figures 5 and 6 show the generated report along with the performance of the chat-bot interface on a few report based and mental health queries, respectively. Furthermore, attempting transfer learning for text classification using BERT and Optuna proved to be a promising approach for the future of mental health classifiers as it gave an accuracy of 80.00% on the initial attempt (Table 1).
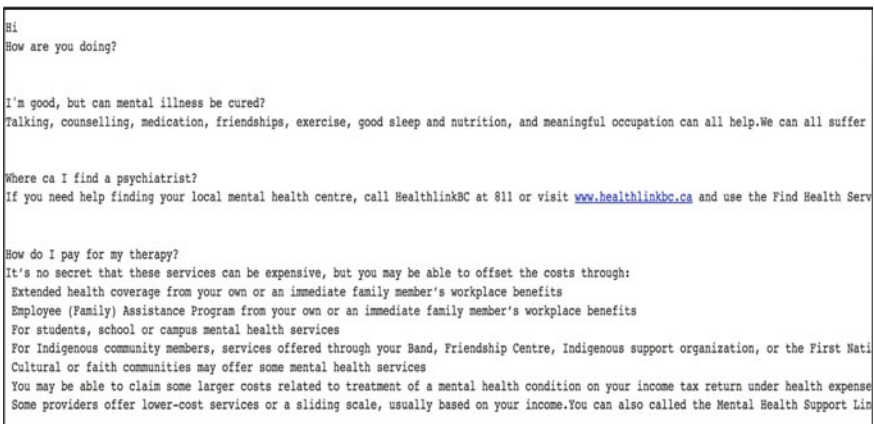


Fig. 6  Bot's response to FAQ's on mental health after summarization

**Table 1** Comparison of classification models

| S. No. | Details | |
|---|---|---|
| | Model | Accuracy (%) |
| 1 | Support Vector Machine (SVM) | 87.87 |
| 2 | Long Short-Term Memory (LSTM) | 89.00 |
| 3 | Transfer learning | 80.00 |

## 6  Conclusion

A declining mental health is an extremely dangerous situation and must be dealt with utmost care and attention. Having a software application to provide you with a detailed report of whether your social media activity indicates declining mental health would be of immense help to the large population that constantly uses such platforms. By utilizing different ML models like Support Vector Machines (SVM) and Long Short-Term Memory (LSTM) model, the proposed framework was able to obtain an accuracy of 88% for the classification. Theoretically, the LSTM models recorded a slightly higher accuracy, but it was found that the Support Vector Machine model worked significantly better on real-time data. Moreover, implementing a transfer learning model using BERT showed promising results on the multi-class classification of tweets and can be explored further.

This proposed approach for an application or plugin that can help keep a user's mental health status in check is only a preliminary framework. It could be extended by enhancing the sentiment analysis of tweets to be able to incorporate sarcasm and emoticon usage into its prediction. Furthermore, the chatbot could be trained to be more context-aware while generating answers.

## References

1. Yatapala KYDHT, Kumara BTGS (2021) Detection of suicide ideation in Twitter using ANN. In: 2021 6th International conference on information technology research (ICITR), pp 1–5. https://doi.org/10.1109/ICITR54349.2021.9657404
2. Tiwari S, Verma A, Garg P, Bansal D (2020) Social media sentiment analysis on Twitter datasets. In: 2020 6th International conference on advanced computing and communication systems (ICACCS), pp 925–927. https://doi.org/10.1109/ICACCS48705.2020.9074208
3. Jain S, Narayan SP, Dewang RK, Bhartiya U, Meena N, Kumar V (2019) A machine learning based depression analysis and suicidal ideation detection system using questionnaires and Twitter. In: 2019 IEEE students conference on engineering and systems (SCES), pp 1–6. https://doi.org/10.1109/SCES46477.2019.8977211
4. Kariya C, Khodke P (2020) Twitter sentiment analysis. In: 2020 International conference for emerging technology (INCET), pp 1–3. https://doi.org/10.1109/INCET49848.2020.9154143
5. Dinakar S, Andhale P, Rege M (2015) Sentiment analysis of social network content. In: 2015 IEEE International conference on information reuse and integration, pp 189–192. https://doi.org/10.1109/IRI.2015.37

6. Amer E, Hazem A, Farouk O, Louca A, Mohamed Y, Ashraf M (2021) A proposed Chatbot framework for COVID-19. In: 2021 International mobile, intelligent, and ubiquitous computing conference (MIUCC), pp 263–268. https://doi.org/10.1109/MIUCC52538.2021.9447652
7. Hwerbi K (2020) An ontology-based chatbot for crises management: use case coronavirus. arXiv:abs/2011.02340
8. Lalwani T, Rathod V (2018) Implementation of a Chatbot system using AI and NLP. Int J Innov Res Comput Sci Technol (IJIRCST) 6(3). ISSN: 2347-5552
9. Das T (2022) Multi-label text classification using transfer learning powered by "Optuna" [Online]. Available: https://www.analyticsvidhya.com/blog/2022/01/multi-label-text-classific ation-using-transfer-learning-powered-by-optuna/. Accessed 10 May 2022
10. Cer D, Yang Y, Kong S, Hua N, Limtiaco N, St. John R, Constant N, Guajardo-Cespedes M, Yuan S, Tar C, Sung Y-H, Strope B, Kurzweil R (2018) Universal sentence encoder. arXiv: 1803.11175