

Social Media Interaction based Mental Health Analysis with a Chat-bot User Interface

Aliyah Kabeer*, Paul John[†], Serena A. Gomez[‡]

*^{†‡} Department of Computer Science and Engineering

PES University, Bangalore, India

{pes2ug19cs031, pes2ug19cs275, serenagomez}@pesu.pes.edu

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 they interact with online may subconsciously have a significant effect on their mental health. This paper presents an 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. To classify the mental state, feature vectors are generated using TF-IDF. Supervised machine learning algorithms like SVM and neural network based models like LSTM are compared for their performance on prediction. Natural Language Processing techniques such as question similarity and extractive summarization are utilized in building the chat-bot framework.

Index Terms—Social Media, Mental Health, Report, Chatbot, Question Similarity, Twitter, Tweets

I. INTRODUCTION

Mental health is a state of emotional and cognitive well-being which affects a person's thought process, feelings and actions whilst coping with the normal stresses of everyday life. In today's world, keeping track of one's 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, we can classify the tweets of a person into various categories such as anxious, stressed, suicidal, lonely, etc. Sentiment Analysis can be performed using various machine and deep learning algorithms such as Recurrent Neural Networks, Support Vector Machines, Logistic Regression, and so on. By fixing a certain threshold value, we can classify each tweet into one of these classes and repeatedly generate a report based on these tweets over fixed periods of time. By generating the report over a

set maximum number of tweets, we can ensure that the report would be balanced.

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. 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, would help 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, we can ensure 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 data-set every time.

The contents of this paper are organized in the following sections: Section II talks about the scope and motivation of this project. In Section III an overview of the related work and previous research pertaining to the concepts used in our paper are reviewed and surveyed. Section IV contains the proposed methodology with subsections for each major module. Section V presents the results, with a comparison of the accuracies of the two classification models used. Section VI and VII contain details of the conclusion and future scope of our work respectively.

II. 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

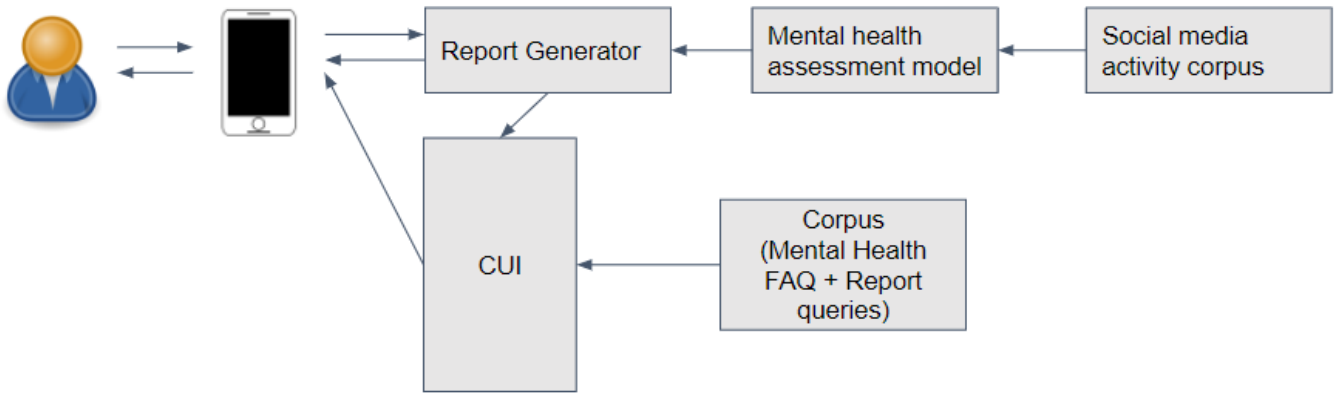


Fig. 1. High Level Architecture.

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, we aim to come up with a virtual assistant plugin that will allow individuals to keep track of their social media usage and assess their mental state along the way.

III. RELATED WORK

Of late, mental health analysis and prevention of the onset of mental health illnesses has been a trending research topic and is 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 has 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. K.Y.D.H.T. Yatapala et al [7] presents a machine learning approach using artificial neural networks to detect suicide ideation or thoughts and identify patterns in suicidal text by analysing 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 TF-IDF worked better and obtained a testing accuracy of 89%. Shikha Tiwari and Anshika Verma [8] 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%. Swati Jain et al [10] also talk about suicidal ideation and depression detection using Twitter data and questionnaires.

They utilised classification algorithms like Random Forest, XGBoost, Support Vector Machine, and Logistic Regression to classify the severity of the detected depression into 5 stages. Although they analyse 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. C. Khariya and P. Khodke in [12] utilised the Twitter API to fetch tweets and classify them into positive, negative or neutral sentiments by using algorithms such as K Nearest Neighbours 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. S. Dinakar, P. Andhale and M. Rege [4] 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 we 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 utilised pattern matching and substituted methodology to make users believe they were talking to a human and give them a safe space to share their feelings. E.Amer et al in [1] 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 Khoulood [3]

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 stopwords 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. Tarun Lalwani and Vasundhara Rathod [9] 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.

IV. PROPOSED METHODOLOGY

In this paper we aim to generate a mental health report based on the users social media activity and provide them with the ability to ask any questions regarding the generated report and also general mental health related questions via a chat user interface. Fig. 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 pre-processed 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 is then used to train the chatbot component of the user interface, which can answer the user's queries post report generation. Fig 2. shows the logic flow followed by the bot for answer generation. The detailed design of each module is as given below:

A. 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.

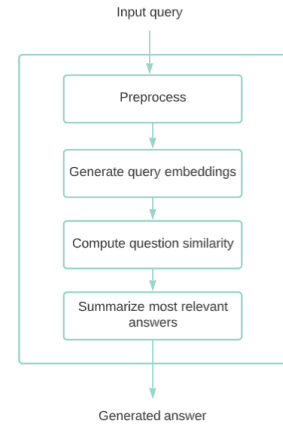


Fig. 2. Logic flow for answer generation.

B. Mental health assessment model

1) *Data Collection*: Twitter being one of the most popular and widely used social media platforms globally, the number of tweets generated every second is unbelievably high. This serves as a continuous source of information for analyzing a user's state of mind over a period of time. The social media activity corpus contains labeled tweets which will be used as the training dataset.

2) *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 stopwords, numbers, punctuation, and special characters are removed. This forms a key part of preprocessing as tweets are generally comprised of various special characters like hashtags, hyper-links, urls and numbers, which might not be relevant to our 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.

3) *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.

4) *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 non-linear 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 labelled 'lonely'. This anomaly could be because there is relatively lesser amount of data available for that particular class in the training dataset.

	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

- *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. LSTM's 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, LSTM's can see the text as an entire sequence and work from there. The LSTM model for the proposed architecture 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 5 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. Fig 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.

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.

C. 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 date of the analysis, along with the Twitter developer credentials. The tweets are then cleaned, preprocessed, vectorized and stored in a dataframe 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.

D. Corpus (Mental health FAQ + Report queries)

The mental health FAQ and report query corpora comprises of 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).



Fig. 5. Report generation followed by the CUI.

E. 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 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 [11] 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 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 is selected to summarize the answer.
- If both the above conditions are not satisfied then the bot utilizes Python’s pre-trained 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,

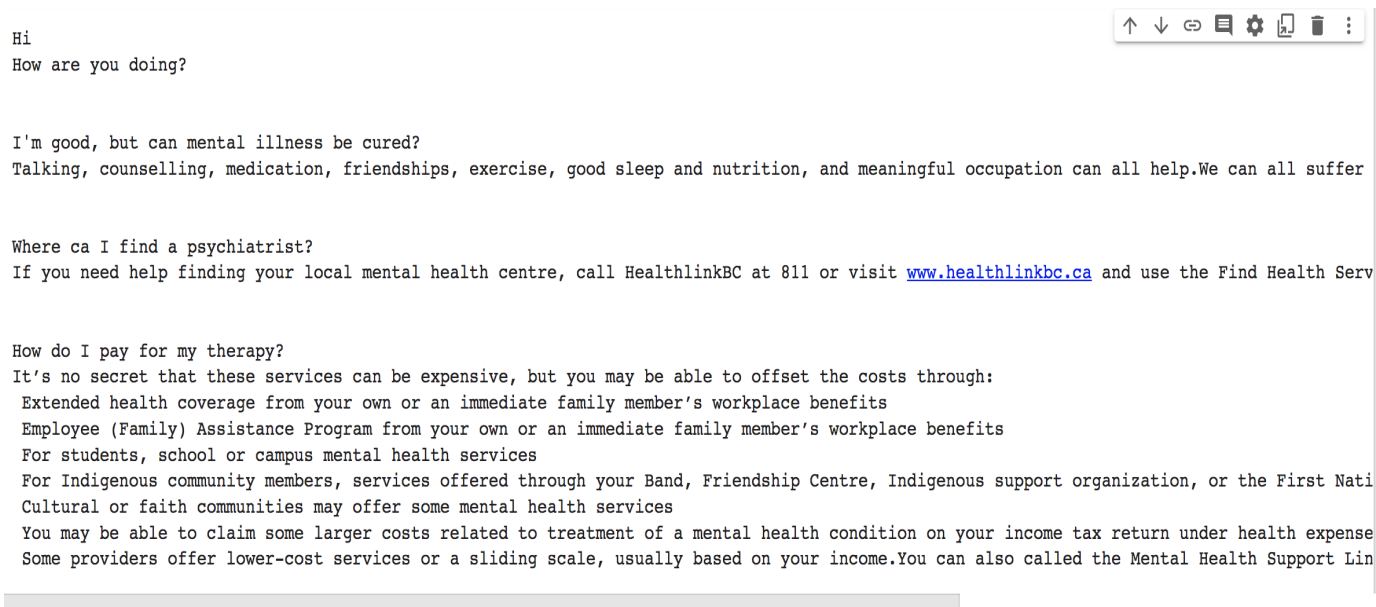


Fig. 6. The bot's response to FAQ's on mental health after summarization.

humor, politics, etc. 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.

V. 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 chatbot interface performed relatively well on domain specific questions like mental health FAQ's and report related queries. Fig 5. and Fig. 6 show the generated report along with the performance of the chat-bot interface on a few report based and mental health queries respectively.

TABLE I
COMPARISON OF CLASSIFICATION MODELS

Sl No.	Details	
	Model	Accuracy
1	Support Vector Machine (SVM)	87.87%
2	Long Short Term Memory (LSTM)	89.00%

VI. CONCLUSION

A declining mental health is an extremely dangerous situation and must be dealt with 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 and Long Short Term Memory model, we were 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 realtime data.

VII. FUTURE SCOPE

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] E. Amer, A. Hazem, O. Farouk, A. Louca, Y. Mohamed and M. Ashraf, "A Proposed Chatbot Framework for COVID-19," 2021 International Mobile, Intelligent, and Ubiquitous Computing Conference (MIUCC), 2021, pp. 263-268, doi: 10.1109/MIUCC52538.2021.9447652.
- [2] Karl Daher, Jacky Casas, Omar Abou Khaled, and Elena Mugellini. 2020. Empathic Chatbot Response for Medical Assistance. Proceedings of the 20th ACM International Conference on Intelligent Virtual Agents. Association for Computing Machinery, New York, NY, USA, Article 15, 1–3. DOI:https://doi.org/10.1145/3383652.3423864
- [3] Hwerbi, Khouloud. "An ontology-based chatbot for crises management: use case coronavirus." ArXiv abs/2011.02340 (2020): n. Pag.
- [4] S. Dinakar, P. Andhale and M. Rege, "Sentiment Analysis of Social Network Content," 2015 IEEE International Conference on Information Reuse and Integration, 2015, pp. 189-192, doi: 10.1109/IRI.2015.37.
- [5] Ayanouz, Soufyane Anouar Abdelhakim, Boudhir Benhmed, Mohammed. (2020). A Smart Chatbot Architecture based NLP and Machine Learning for Health Care Assistance. 10.1145/3386723.3387897.
- [6] Suresh Raj, Vivek. (2021). Performance of Seq2Seq learning Chatbot with Attention layer in Encoder decoder model. 10.13140/RG.2.2.33355.92961.
- [7] K. Y. D. H. T. Yatapala and B. T. G. S. Kumara, "Detection of Suicide Ideation in Twitter using ANN," 2021 6th International Conference on Information Technology Research (ICITR), 2021, pp. 1-5, doi: 10.1109/ICITR54349.2021.9657404.
- [8] S. Tiwari, A. Verma, P. Garg and D. Bansal, "Social Media Sentiment Analysis On Twitter Datasets," 2020 6th International Conference on Advanced Computing and Communication Systems (ICACCS), 2020, pp. 925-927, doi: 10.1109/ICACCS48705.2020.9074208.

- [9] Tarun Lalwani, Vasundhara Rathod, Implementation of a Chatbot System using AI and NLP, International Journal of Innovative Research in Computer Science Technology (IJIRCST) , Vol-6, Issue 3., May 2018, ISSN 2347 - 5552
- [10] S. Jain, S. P. Narayan, R. K. Dewang, U. Bhartiya, N. Meena and V. Kumar, "A Machine Learning based Depression Analysis and Suicidal Ideation Detection System using Questionnaires and Twitter," 2019 IEEE Students Conference on Engineering and Systems (SCES), 2019, pp. 1-6, doi: 10.1109/SCES46477.2019.8977211.
- [11] Cer, Daniel and Yang, Yinfei and Kong, Sheng-yi and Hua, Nan and Limtiaco, Nicole and John, Rhomni St. and Constant, Noah and Guajardo-Cespedes, Mario and Yuan, Steve and Tar, Chris and Sung, Yun-Hsuan and Strophe, Brian and Kurzweil, Ray, "Universal Sentence Encoder", 2018 arXiv:1803.11175.
- [12] C. Kariya and P. Khodke, "Twitter Sentiment Analysis," 2020 International Conference for Emerging Technology (INCET), 2020, pp. 1-3, doi: 10.1109/INCET49848.2020.9154143.