

# A Deep Learning Framework Using Data Augmentation for Accuracy Improvement to Analyze Users Posts on Social Media to Find Signs of Mental Illness

Vaibhav Sharma<sup>1</sup>, Parul Goyal<sup>2\*</sup>

<sup>1</sup>School of CA & IT, Shri Guru Ram Rai University, Dehradun, Uttarakhand-248001, India  
Email:- vsdeveloper10@gmail.com, Orcid-0000-0002-1404-2012

<sup>2\*</sup>School of CA & IT, Shri Guru Ram Rai University, Dehradun, Uttarakhand-248001, India  
Email:- profdrparul@gmail.com, Orcid-0000-0001-9729-1155

\*Corresponding Author-Parul Goyal

<sup>1</sup>School of CA & IT, Shri Guru Ram Rai University, Dehradun, Uttarakhand-248001, India  
Email:- profdrparul@gmail.com, Orcid-0000-0001-9729-1155

## Abstract-

Through their posts on social media, users frequently express their feelings. A deep learning model was developed for this study to determine a user's mental state based on the data they posted. We gathered articles for this purpose from Reddit forums dedicated to mental health. Our suggested model may disorder, anxiety, depression and Schizophrenia by examining and learning posting information published by users. Based on their posts, we think our algorithm can help identify people who could be experiencing mental illness. The consequences of this model, which may be used in combination to other methods to track the mental health of individuals who use internet extensively, are also discussed in this paper.

**Keywords-** Mental health, Mental illness, Anxiety, Depression, Schizophrenia etc.

## 1. Introduction

Users frequently convey their emotions on social media[1]. Users are frequently more prone to express their mental illnesses or difficulties through various social media or online social health networks[2]. By interacting with people who have comparable symptoms, these online health forums can provide as a network for expressing sympathy. In an effort to self-diagnose, users frequently seek for medical information on social media that relates to their symptoms[3]. In line with this development, a number of academics have examined user-generated social media content in order to observe users' emotional states or mental illnesses, such as sadness, schizophrenia depression and anxiety[5][6]. A new study obtained tweets from individuals who were said to have depressed illnesses and used the Linguistic Inquiry and Word Count (LIWC)[7] to analyze the tweets' linguistic and emotional content monitored the users' shifts in social activity on Twitter, advanced psychometric techniques were employed in a different investigation to contrast the rates of depressive symptoms between the pre- and post-natal periods in order to predict users' postpartum depression based on their Facebook posts and comments. Additionally, Reece et al.[8] used image data to find depressed users on social networking sites. Both face detection and colorimetric analysis were used after collecting photographs that users had uploaded to Instagram. Prior research that gathered user

data from Reddit users and demonstrated the effectiveness of using N-gram phrase modeling, vector embedding methods, and topic monitoring of user postings, it is possible to spot prospective sufferers of anxiety and depression among users.

It has been demonstrated that data from social media can be useful in observing or detecting individuals' moods or possible mental state health issues in a number of earlier research. This work takes it a step further; we seek to create a deep learning model that can recognize a user's mental problem, specially depression, anxiety and, schizophrenia by gathering various mental-health-related data from social media. In order to achieve this, we gathered user postings from Reddit, a famous social media site with a variety of communities (or so-called "subreddits") devoted to mental health, including depression, bipolar, and schizophrenia.

We received data from the three subreddits depression, Since our objective is really to identify how a user seems to have a mental disorder, including depression and anxiety, these two conditions are considered. Note that we used the mental-health-related Reddit, subreddits that were previously discovered. More specifically, a statistical method like a semi-supervised method and an expert assessment procedure were used to identify 3 prominent out of many subreddits as being connected to mental health. Each subreddit that has been found is linked to a particular mental illness.

Stream	% of users	% of posts	Description
Depression	54.2%	52.9%	Anyone suffering from depression or another mental disease should receive peer assistance.
Anxiety	19.8%	17.7%	For those who have any type of anxiety problem, there is discussion and help available.
Schizophrenia	2.1%	3.6%	Hello! The purpose of this community is to explore schizophrenia spectrum diseases and relevant topics. You are free to publish, talk, or simply lurk. We are all here for one another, so there is no place for judgement. Please refrain from self-diagnosing, advising specific medical procedures, or diagnosing others.

**Table1.**An overview of the data gathered from Reddit.

We looked at whether certain posts made by the user may be categorised as relevant sorts of mental disorders by gathering and examining the user's posts that were submitted in various Reddit subreddits pertaining to mental health. People with certain mental disorders may not be aware of their most accurate diagnosis; for instance, people with bipolar disorder may find it difficult to differentiate between the disorder and depression because of the similar symptoms, or even the initial diagnosis of bipolar disorder may be challenging[9]. We supposed most users seek to browse for generic terms related to by presenting early-stage general concerns about psychological health on social media, such as "mental wellbeing," "mental illness," or "mental stability," it appears as though people are seeking out for remedies. As a result, a lot of users prefer to engage in conversation in one of Reddit's general wellbeing channels (for example, mentalhealth) at first, but frequently refuse to recognise their true issues. As a result, we try to identify users' probable mental illnesses from their social media postings. The mentioned research question will be examined in the study.

**Research question:** Can we determine whether a user's social media posts are related to mental illnesses or mental health?

## 2. Methods for Study

**2.a) Collection of Data-**We gathered post information from the following three sub-reddits dedicated to mental health, of which sadness, anxiety, and schizophrenia are each believed to be connected to a certain disorder. In order to research postings that offered general health-related information, we also collected post data from the most popular health-related sub-reddit, mentalhealth [10]. Every user ID from every subreddit whose postings included at least one mental health-related topic was collected. With the help of the PushshiftAP, we gathered post titles in addition to user IDs[11].

It should be noted that all user data was anonymized; as a result, no personally identifying information was included. The details of the information collected is summarised in Table 1.

**2.b) Data Pre-Processing Steps-**It shows how to pre-process the post-data that was collected. Each title's matching post was combined after the data had been gathered. For each post, we eliminated extra punctuation

and blank spaces. Next, we tokenized user postings using Python's natural language toolkit (NLTK) and filtered words that were often used (stop words). On the tokenized words, A tool was created to create a set of guidelines for analysing the source or meaning of words that applied in order to reduce the word corpus and convert each word to its root meaning.

**2.c) Classification Models-**We created three binary classification models, each of which assigns a user's individual post to each of the subreddits for depression, anxiety and schizophrenia. Our hypothesis is that a member with a particular mental illness posts a message on the relevant subreddit that addresses the issue. If a user has many mental health issues, such as depression and anxiety, for example, he or she may post in various subreddits. The categorization model, however, may be affected by noisy data if the model is trained using posts from users who have various symptoms, as in a previous study. Therefore, in order to enhance performance, we created six distinct binary classification models for each symptom. We were successful in effectively predicting a user's potential by developing three unique models for every mental illness, each of which incorporates information from individuals who have just one type of mental illness. For instance, we defined the depression class as the posts submitted by users who only contribute content within the depression category in order to construct a model for identifying depression. The opposing class is known as the non-depression category. We used the synthetic minority over-sampling method (SMOTE) approach to correct a problem with class imbalance in the data that was collected.

Our dataset was split into training (79%) and testing (21%) sets. Convolutional neural network (CNN) was then used. Additionally, during the learning phase, we disregarded posts from people who published in numerous sub-reddits. We changed the words in the training set to numerical representations to quantitatively represent each post (Fig. 1). Using the word2vec API of the Python Package, Gensim, we employed word-embedding techniques from the pre-processed texts in the case of the CNN classifier. The continuous bag-of-words representation (CBOW) models were used to pre-train the word vectors using the training dataset gathered for the current investigation, with a window size of 5. It should be noted that by employing the word2vec model for describing each post for every topic, a linguistic

style utilised by authors of posts in a specific subreddit may be learned.

Fig. 1 displays a generalised representation of the suggested CNN-based model. The layers that make up the model architecture are ordered in a particular order and consist of a dense layer, an embedding layer, a fully connected layer, a layer using maximum pooling, and the output. Fig 1 illustrates the training process for a post using the given model. A pre-trained word2vec initialises the strength of embedding layer, the initial layer in the model, which reflects the word representations of a pre-processed post with 20 dimensions. Second, a convolutional layer with word vector input comprises 128 filters with a five-filter-size increment between each filter. In order to avoid over-fitting problems, we also added a dropout rate of 25%. The max-pooling layer has a 128-dimensional dimension and uses the

maximum values found in the CNN filters. The result is the sigmoid activation function's probability of the categorization, which runs from 0 to 1. It is obtained by passing the output of the max-pooling layer through two dense layers that are fully connected.

With a learning rate of 0.001, we trained the neural network utilising the Adam optimizer and binary cross-entropy loss function. The batch size was set at 64 and our model was trained across 50 epochs.

### 3. Results

To verify the effectiveness of the models, four evaluation measures were used which are as follows accuracy, precision recall, and F1-score .  $A$ (represent true positive),  $B$ (false negative),  $C$ (true negative), and  $D$ (falsepositive).

$$\text{accuracy} = \frac{TP + TN}{TP + FN + TN + FP} \quad (1)$$

$$\text{recall} = \frac{TP}{TP + FN} \quad (3)$$

$$\text{precision} = \frac{TP}{TP + FP} \quad (2)$$

$$\text{f1-score} = \frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}} \quad (4)$$

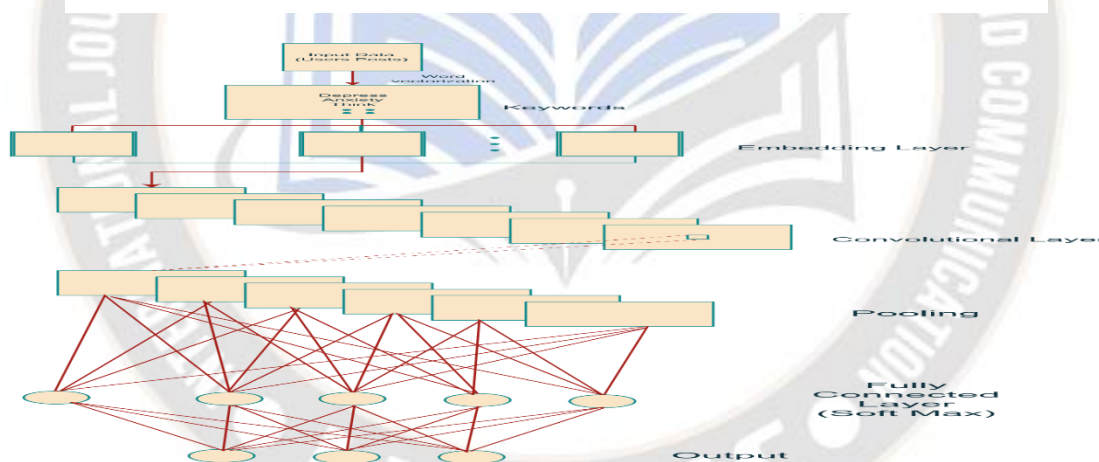


Fig 1. An architecture of the proposed model.

The results of the three binary classification models are collected in Table 2 and Fig 2 and Fig 3 shows accuracy and comparisons with previous model . Across all subreddits, proposed Deep learning models performed more accurately than previous model overall. Precision (92.46%), recall (86.17 %), and F1-score (80.72) performance ratings for the depression class were at their highest on one of the most evenly distributed subreddits, depression. Additional subreddits—Anxiety, —also demonstrated high accuracy with CNN models, with scores of 82.34%. However, these subreddits' F1-scores for identifying mental illnesses ranged from 50-65%, which is considerably lower than the scores from the class-balanced channels. In conclusion, our suggested approach can reliably identify potential users who

might be suffering from psychiatric illnesses. We think that by gathering more data, the issue of imbalanced data may be solved, leading to improved performance.

### 4. Discussion

Potential patients of mental disorders can be helped by identifying problems with mental disease early on and offering suitable solutions[12]. We developed a deep learning model utilising natural language processing to identify indicators of mental illness by gathering and examining data from social media platforms pertaining to mental health individuals who may be suffering from mental illness based on their posts subreddits on Reddit that concentrate on issues relating to mental disorders. We think

that our approach can use in a new era of study where online social media can serve as a useful tool for spotting probable mental illnesses based on users' particular posts[13].The majority of those who may have mental illnesses, however,

nevertheless experience social blind spots and do not receive the proper care due to a variety of factors, such as finding it difficult to disclose their condition to someone in person or having trouble physically getting to clinics[12].

Channel	CNN model				Proposed model			
	Precision	Recall	F1-Score	Accuracy(%)	Precision	Recall	F1-Score	Accuracy(%)
Depression	90.40	69.53	80.64	74.75	92.46	86.17	80.72	83.45
Anxiety	86.67	40.65	55.53	78.65	90.63	55.14	64.30	82.34
Schizophrenia	82.54	25.32	39.65	93.66	84.62	26.42	40.55	94.72

Table 2: CNN evaluation and the proposed model.

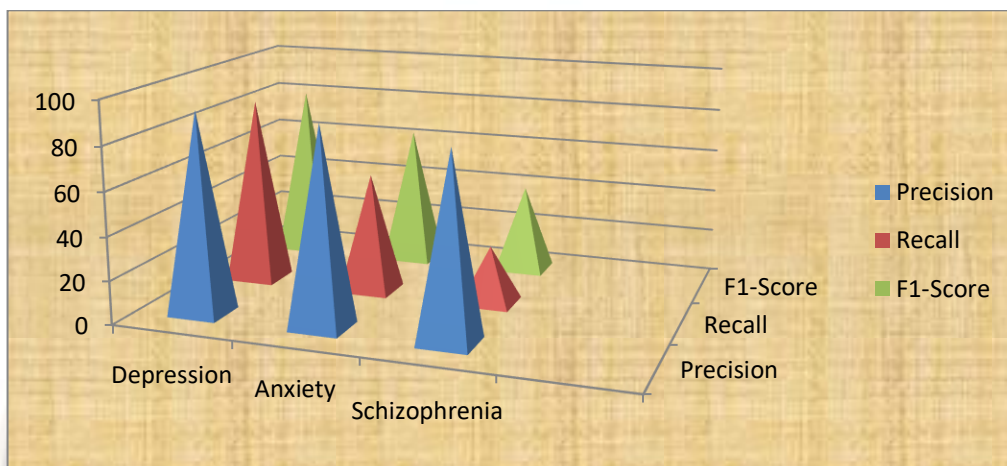


Fig 2- Proposed Model

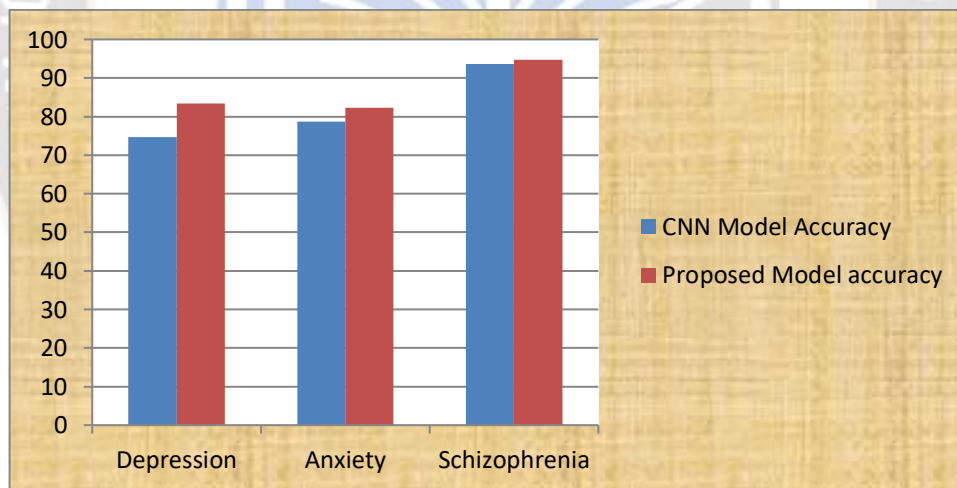


Fig-3-Accuracy comparison CNN vs Proposed Model

The following implications are provided in light of the lessons discovered. First, it is possible to identify potential mental diseases in users by reading their postings when deep learning algorithms are combined with the right natural language processing techniques. The methods utilised in this study can be used to warn people who may be suffering from particular mental diseases before they attend counselling centres by making use of readily available social media data[15]. Second, this study offers compelling data that suggests it may be possible to use online tools to support those in need of psychiatric assistance. For illustrate, before presenting the chances of

each mental illness indicated by our certified models based on the user's writings, internet site service providers may first ask a user's consent to access their account. The most recent study suggests that recognising mental disorder on social networks sites may develop into a popular area of study in the future. The findings of this study indicate that social media sites could help create a space where people who are suffering from mental illness can talk about their experiences can engage with one another. This study does have certain restrictions, though. The socio-demographic and regional variables that may have an impact on the classification models were not taken into

account in the current investigation. Future studies that aim to enhance the deep learning models' precision or quality can take these parameters into account. Additionally, we gathered information from Reddit, a public social network that may express user sentiments differently than the user's personal social network feed. We did not do further validation procedures for our model, as previously stated. Using a different independent dataset, which will require more research. Online social media has despite the fact that blog posts there can't be as explicit as those on users' personal pages, which may imply they have indeed been diagnosed with diagnostic psychological disorders, users share their symptoms fairly accurately under the semi-anonymity system, giving the site the potential to be utilized to recognise people with mental illness. In order to directly diagnose the symptom and provide the anticipated probability for each symptom, we also trained our model on a particular mental state. As a result, we were unable to quantify other mental illness status adequately; this task will be left for further study[16].

Our multiple binary classification algorithms could be used in future study to identify additional real-world mental disorders. Additionally, we want to use online articles to test our proposed approach, by users who might be suffering from an undiagnosed mental illness on Facebook or Twitter. Additionally, a RNN-based user-level mental disease detection model can be created with the aid of a time-series user-level analysis that monitors a user's longitudinal behaviour pattern.

## References

1. Al-Saggaf, Y., & Nielsen, S. Self-disclosure on facebook among female users and its relationship to feelings of loneliness. *Comput. Hum. Behav.* 36, 460–468 (2014)
2. Shen, J.H., & Rudzicz, F. Detecting anxiety through reddit. In *Proceedings of the Fourth Workshop On Computational Linguistics and Clinical Psychology-From Linguistic Signal to Clinical Reality* (2017).
3. Yoo, M., Lee, S., & Ha, T. Semantic network analysis for understanding user experiences of bipolar and depressive disorders on reddit. *Inf. Process. Manag.* 56, 1565–1575 (2019).
4. Park, A., & Conway, M. Harnessing reddit to understand the written-communication challenges experienced by individuals with mental health disorders: analysis of texts from mental health communities. *J. Med. Internet Res.* 20, e121 (2018).
5. Ernala, S.K., Rizvi, A.F., Birnbaum, M.L., Kane, J.M., & DeChoudhury, M. Linguistic markers indicating the therapeutic outcomes of social media disclosures of schizophrenia. *Proc. ACM Hum.-Comput. Interact.* 1, 43 (2017).
6. Gkotsis, G. et al. Characterisation of mental health conditions in social media using informed deep learning. *Sci. Rep.* 7, 45141 (2017).
7. Pennebaker, J. W., Boyd, R. L., Jordan, K., & Blackburn, K. Te development and psychometric properties of liwc2015. Technical Report (The University of Texas at Austin, Austin, 2015).
8. Reece, A.G., & Danforth, C.M. Instagram photo reveal predictive markers of depression. *EPJ DataSci.* 6, 15 (2017).
9. Huang, Y.-H., Wei, L. H., & Chen, Y.S. Detection of the prodromal phase of bipolar disorder from psychological and phonological aspects in social media. *arXiv:1712.09183* (2017).
10. Gaur, M. et al. "Let me tell you about your mental health!" contextualized classification of reddit posts to dsm-5 for web-based intervention. In *Proceedings of the 27<sup>th</sup> ACM International Conference on Information and Knowledge Management*, 753–762 (2018).
11. Pushshif.io. Pushshif.io: Learn about big data and social media ingest and analysis. <https://pushshif.io/> (2019)
12. Hunt, J., & Eisenberg, D. Mental health problems and help-seeking behavior among college students. *J. Adolesc. Health* 46, 3–10 (2010)
13. Wicks, P. et al. Perceived benefits of sharing health data between people with epilepsy on an online platform. *Epilepsy Behav.* 23, 16–23 (2012).
14. Antoniou, Grigoris, Emmanuel Papadakis, and George Baryannis. "Mental Health Diagnosis: A Case for Explainable Artificial Intelligence." *International Journal on Artificial Intelligence Tools* 31.03 (2022): 2241003.
15. Balasubramanian, S., Devarajan, H. R., Raparathi, M., Dodda, S. B., Maruthi, S., & Adnyana, I. M. D. M. (2023). Ethical Considerations in AI-assisted Decision Making for End-of-Life Care in Healthcare. *PowerTech Journal*, 47(4), 168. <https://doi.org/10.52783/pst.168>
16. Singh, Palak, et al. "Artificial Intelligence based Early Detection and Timely Diagnosis of Mental Illness-A Review." 2022 International Mobile and Embedded Technology Conference (MECON). IEEE, 2022.
17. Zaved, M., Ahmed, I., Sinha, N., & Phadikar, S.. "Automated Feature Extraction on AsMap for Emotion Classification using EEG". <https://doi.org/10.48550/arxiv.2201.12055>, 2022.