

Depression Intensity Estimation via Social Media: A Deep Learning Approach

Shreya Ghosh^{id}, Graduate Student Member, IEEE, and Tarique Anwar^{id}

Abstract—Depression has become a big problem in our society today. It is also a major reason for suicide, especially among teenagers. In the current outbreak of coronavirus disease (COVID-19), the affected countries have recommended social distancing and lockdown measures. Resulting in interpersonal isolation, these measures have raised serious concerns for mental health and depression. Generally, clinical psychologists diagnose depressed people via face-to-face interviews following the clinical depression criteria. However, often patients tend to not consult doctors in their early stages of depression. Nowadays, people are increasingly using social media to express their moods. In this article, we aim to predict depressed users as well as estimate their depression intensity via leveraging social media (Twitter) data, in order to aid in raising an alarm. We model this problem as a supervised learning task. We start with weakly labeling the Twitter data in a self-supervised manner. A rich set of features, including emotional, topical, behavioral, user level, and depression-related n -gram features, are extracted to represent each user. Using these features, we train a small long short-term memory (LSTM) network using Swish as an activation function, to predict the depression intensities. We perform extensive experiments to demonstrate the efficacy of our method. We outperform the baseline models for depression intensity estimation by achieving the lowest mean squared error of 1.42 and also outperform the existing state-of-the-art binary classification method by more than 2% of accuracy. We found that the depressed users frequently use negative words such as stress and sad, mostly post during late nights, highly use personal pronouns and sometimes also share personal events.

Index Terms—COVID-19, deep learning, depression intensity estimation, mental health, social media mining.

I. INTRODUCTION

DEPRESSION is one of the most overlooked causes of suffering and death globally, especially among young adults. It is important to understand and realize the far-reaching negative impacts of this silent killer [1]. According to WHO [2], more than 264 million people are affected worldwide and it is increasing day by day. Depression is totally different from the usual mood fluctuations and ephemeral emotional responses that we face in our everyday life [3].

Manuscript received May 12, 2020; revised April 24, 2021; accepted May 12, 2021. (Corresponding author: Tarique Anwar.)

Shreya Ghosh was with the Department of Computer Science and Engineering, IIT Ropar, Rupnagar 140001, India. She is now with the Faculty of Information Technology, Monash University, Melbourne, VIC 3168, Australia (e-mail: shreya.ghosh@monash.edu).

Tarique Anwar is with the Department of Computing, Macquarie University, Sydney, NSW 2109, Australia, and also with CSIRO's Data61, Eveleigh, NSW 2122, Australia (e-mail: tarique.anwar@mq.edu.au).

Digital Object Identifier 10.1109/TCSS.2021.3084154

It may become a serious health issue when it lasts more than two weeks with moderate or severe intensity. It can affect a person's mind as well as physical health. As a result, a depressed person functions poorly at the workplace and misbehaves with family and close ones. It can even lead to suicide if a depressed person is not receiving proper treatment. Each year approximately 800 000 people die due to suicide. It is the second leading cause of death among teenagers. The situation is worst in countries like India, China, and the USA, as compared to the global scenario. India is said to be the most depressed country in the world. According to WHO [4], India, China, and the USA are the worst victims of anxiety, schizophrenia, and bipolar disorder. Despite having proper effective treatment for depression, not many people among the affected ones are able to receive such treatment for various reasons. In many counties, only around 10% approach for treatment. The main reason behind these poor statistics is the social stigma associated with mental disorders. Other barriers associated with effective care include an inaccurate assessment, lack of resources, and lack of trained healthcare providers. Even more than 70% of people would not consult a psychologist in their early stage of depression, which further leads to the deterioration of mental health.

A. COVID-19 and Depression

The current outbreak of coronavirus disease (COVID-19) is affecting us not just physically, but also psychologically. It is a particular and rare situation that is impacting in varying ways on an international scale, including an alteration in our overall lifestyle. The measures of social distancing and *lockdown* implemented in several affected countries have led to interpersonal isolation and extreme changes in our daily lives. Other consequences on the lives of people include loss of employment, financial insecurity, domestic violence, and cybercrimes. These conditions altogether have raised serious concerns of stress, anxiety, and depression [5]. With an overburdened work under stressful circumstances, the health care workers too are heavily prone to depression, anxiety, and insomnia [6]. Pappa *et al.* [6] stressed the need to establish ways to mitigate mental health risks of health care workers under the current pandemic conditions. Psychologists and social scientists are contemplating that COVID-19 could lead to an epidemic of clinical depression and our health care system is way behind for proper treatment [7]. Furthermore, Dresden [8] reports that, just as the COVID-19 can be much more dangerous with

preexisting medical conditions, the psychological effects of the pandemic can be much worse for people with preexisting depression. In the presence of this highly infectious disease and a mental health care system with limited services to deal with large mental health issues, the demand for remote therapy is rapidly rising [9]. A text and video chat therapy service called Talkspace has already seen a 65% increase in customers since mid-February 2020. Brightside is a mobile app that offers treatment and medication for anxiety and depression. It has seen a 50% bump in new users since the start of the quarter. A digital therapeutics company called Big Health is releasing cognitive and behavioral technique-based programs to combat poor sleep and anxiety for free. More than 50 companies have signed up or expanded their use of programming including big employers such as Nike, Target, and supermarket chain HEB. With this high demand for depression treatment with limited resources under the extreme circumstances of COVID-19, it urges for rapid development of technology-based support in the treatment process.

B. Social Media as a Tool for Depression Detection

Nowadays, people extensively use social media platforms like Facebook, Twitter, Weibo, and WhatsApp. These social media sites have become a platform for users to express their views, ongoing moods, and emotions and share with their family, friends, and other related people. The behaviors of posting and sharing content on social media reflect the users' daily lives and their mental state. In the past few decades, researchers have widely used social media for detecting the mental wellness of a user [10]–[12] and various other applications [13]–[16]. Park *et al.* [17] conducted an interview-based and questionnaire-based user study to analyze the behaviors and the language patterns of depressed users on Twitter. These interview-based methods are expensive and time-consuming. Moreover, it is often difficult to get sufficient data in this manner to guarantee the robustness and generalized character of the model. These interview questionnaires are based on depression behaviors, as established in the previous related research [18]. The symptoms of depression also evolve with the advancement of technologies, especially after the advent of online social media. These evolving symptoms may not always be totally covered in the previous depression literature. To bridge the gap, some studies [19], [20] recently analyzed the social media behavior-based model to detect depression. They collected a Twitter-based depression dataset and developed a method to leverage cross-domain depression information. However, according to the existing literature [18], [21], depression is a disorder that may exist in people at different levels or intensities. At an early stage, this intensity is generally low and its severity is driven by the situation being experienced over time. The affected person generally would require treatment according to the severity of depression. Hence, just detecting depression in a person is not always helpful for recommending the proper treatment by a domain expert. Instead, it is very important to determine the severity of depression for a proper understanding of the case, as illustrated

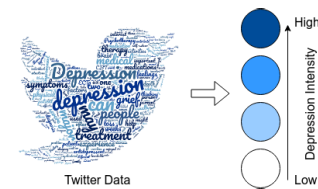


Fig. 1. Depression intensity analysis from social media on BDI-II scale.

in Fig. 1. It would help early-stage depression patients with low intensities to take preventive measures.

C. Our Contributions

In this article, we propose a deep learning-based method to estimate the intensity of depression by using the contents shared by the user on social media platforms. To the best of our knowledge, this is the first such work on depression intensity estimation. To start with, we relabel the depression dataset shared by Shen *et al.* [19] in a self-supervised manner into different intensity categories based on the textual compound polarity and latent semantic analysis (LSA). Considering the different properties of depression established in the literature, we design a total of 527 features of five different types to describe each user, which includes emotional, event-triggered, behavioral, user-level, and depression-related features. Using the extracted features, we train a shallow long short-term memory (LSTM) network to predict the depression intensities. We compare our experimental results with different other models to evaluate our intensity estimation. Furthermore, we also compare our method (adapted to binary classification) with existing binary classification methods and outperform them. In summary, we make the following main contributions in this article.

- 1) We present a comprehensive literature review of this interdisciplinary domain of mental health and social media mining for depression detection.
- 2) We develop a dense semisupervised labeling technique, to relabel the existing sparsely labeled Twitter data, with the intensity of depression.
- 3) We propose a method, consisting of extraction of five different types of relevant features for each user from their social media data and training an LSTM network, to predict the intensity of depression of the users.
- 4) We perform extensive experiments to establish the efficacy of the proposed method. We outperform the other comparable models for intensity estimation and also outperform the existing methods for binary classification.

D. Paper Organization

The remainder of this article is organized as follows. Section II provides a comprehensive literature, followed by some fundamental concepts and the problem statement in Section III. We present the proposed method in Section IV and experimental results in Section V. Finally, we conclude this article in Section VI.

II. LITERATURE REVIEW

The detection of depression from social data is an interdisciplinary domain of depression, mental health, and social

media mining. In this section, we discuss the existing literature and recent advancements in depression detection, starting from traditional methods in the early era to the present day of leveraging social media data.

A. Traditional Depression Detection

Realizing the seriousness of the problem of depression, its study started much earlier than the Internet era [18], [22]. These studies were mainly based on user-dependent questionnaire surveys. The *depression inventory* of Beck *et al.* [18] comprises 21 questions, which mainly judges the users' mental state. Similarly, the CES-D Scale [22] is also based on 20 questions. These questions are mainly on mental conditions, like users' guilty feelings and sleeping patterns. The questionnaires either have multiple choice type questions with variable scores or ask for users' feedback along with the degree of their situations. The depression level is further diagnosed based on the total score. Furthering the detection of depression, its intensity estimation is also being studied since long ago. Beck Depression Index-II (BDI-II) [21] is a standard in "broad ranges," which maps the questionnaire-based depression scores into the following categories: 1) 0–13: indicates no or minimal depression; 2) 14–19: indicates mild depression; 3) 20–28: indicates moderate depression; and 4) 29–63: indicates severe depression. On the contrary, Whooley and Owen [23] came up with nine kinds of depressive indicators as a standard criterion for depression diagnosis of their diagnostic and statistical manual of mental disorders (DSM). Clinical psychologists generally follow this method to diagnose depression over a period of time before the final decision. These methods are applied to real-world depression cases for several years. Also to be noted is that the DSM keeps evolving with time. It has evolved from the DSM-IV version to the DSM-V version (2013) in a span of 12 years, where the new social media-based behaviors and symptoms may not have been included yet.

B. Social Media and Depression Detection

As social media has become a popularly accepted platform to express one's inner-self, researchers have recently started leveraging user-generated data on such platforms for depression detection and mental health care. Park *et al.* [17] first analyzed the language used on Twitter for expressing depressive moods. They further investigated [24] on Twitter users' depressive attitudes and behaviors. For their study, they conducted face-to-face interviews and measure the correlation between those interviews and Twitter data. De Choudhury *et al.* [25] also explored social media data for "major depressive disorders" prediction. They leverage several behavioral patterns of depressed users from this data and correlate them with psychological studies. Resnik *et al.* [26] studied the topic models to analyze the linguistic signal patterns corresponding to social engagement, poor ego control, relationship, emotional distress, and so on for depression prediction. Similarly, Tsugawa *et al.* [27] further investigated Twitter posting behavior to predict depression for Japanese-speaking users. Xue *et al.* [28], [29] proposed models for stress detection using social media data. They first detected psychological pressures via analysis of teenagers' tweets. Furthermore, they

assisted them to get relief from their stress through microblogs. A recent work by Xu and Zhang [30] attempted to explain how social media users disclose their moods from the perspective of depression analysis. In another recent work, Shen *et al.* [19] constructed a well-labeled depression dataset having depression and nondepression labels on Twitter. They extracted six depression-related feature groups, which cover both the clinical depression criteria and online behaviors on social media. To achieve the abovementioned task, they proposed a multimodal depressive dictionary learning model that detects depressed users on Twitter. Shen *et al.* [20] further studied an interesting problem of transferring depression labels from a source domain to a target domain. Their work transfers from Twitter as the source domain to Weibo as the target domain. Sadeque *et al.* [31] recently conducted a study to measure the latency of depression detection from social media. It measures how accurately the model can predict depression from Reddit data in terms of early risk detection error (ERDE), risk window-based method, and latency. From the above literature, it is evident that one can use social media behavioral patterns for depression analysis. Similarly, Trotzek *et al.* [32] proposed a deep learning-based method for early detection of depression in a social media platform. They evaluated different word embeddings via CNN and compared the results with user-level linguistic metadata-based classification. There exists several works on influential community identification and influence maximization [33]–[35]. They can help to analyze the role of communities in the propagation of depression across the social network.

C. Event Triggered Mental Health Analytics

People sometimes get affected mentally due to some negative experience in their real life. Such an event acts as a trigger for mental problems, which further leads to depression. Similar to the problem of event detection in general, *personal event detection* from social media data has also become an interesting direction of research. Li and Cardie [36] assembled the users' life events based on their Twitter profiles. Furthermore, they proposed an unsupervised framework to create a list of personal events. Similarly, in another work, Li *et al.* [37] trained a supervised classifier that takes manually defined textual features as input. These features classify users' tweets into predefined life event categories. Some works focus on the stress getting generated on social media. Lin *et al.* [11] worked on identifying the subjects or events that trigger the stress. Based on this hypothesis, they collected a dataset from Twitter and label them properly based on the stress events.

D. Emotion from Social Texts

Emotions are closely related to depression and mental health [38]. The social media trends provide a good opportunity to judge the user emotion from their textual data. Users often share a large volume of posts and tweets, which describe their ongoing emotional state. To interpret the syntactic and semantic meaning of a given textual post, the data are converted to a vector format and mined. There are various methods to create these representations. Keyword or lexicon-based

to estimate the depression intensity of each user on a [0-1] scale that can be easily mapped to a level of depression. We use BDI-II [21], where the corresponding [0-1] scale is subdivided into different ranges as no or minimal depression (0 corresponds to no depression), mild depression, moderate depression, and severe depression (1 corresponds to most-severe depression).

IV. PROPOSED METHOD

The problem of estimating the intensity of depression of a person is challenging, especially because of the noisy and unstructured nature of the social media data. We list some of the main challenges below.

- C.1 There is no publicly available large-scale benchmark dataset for depression intensity analysis.
- C.2 Users' behaviors on social media are heterogeneous. It is difficult to characterize the users from discriminant perspectives and capture the relation across different modalities.
- C.3 Although users' behaviors are rich and diverse, only a few are symptoms of depression. This makes the depression-oriented features sparse on social media and difficult to be captured.

In order to develop our method, we need to address the aforementioned challenges. We address C.1 by developing a dense labeling technique in Section IV-A, to relabel the existing sparsely labeled dataset, with the intensity of depression. We address C.2 and C.3 by developing a deep learning-based method in Sections IV-B–IV-D, in which we preprocess the social data, extract a rich set of features, and train an LSTM network to predict the depression intensity. C.2 is addressed by defining our features corresponding to users, which makes the method independent of their heterogeneous nature. C.3 is addressed by considering a rich set of features having *indirect* relations with depression. Fig. 3 shows the overall pipeline of the proposed method. It starts with dataset preprocessing and relabeling into different levels of depression by computing a depression score for each user. Then a rich set of features are extracted and an LSTM network is trained to predict the final depression intensity of users.

A. Dataset Description and Relabeling

We use the dataset curated by Shen *et al.* [19], consisting of a total of 6562 users, out of which 1402 are labeled as *depressed* and 5160 as *nondepressed* and a total of 4245747 tweets, out of which 292564 are of depressed users and 3953183 are of nondepressed users. They collected this dataset via web-crawling the profile information of the Twitter users along with their timeline. They also searched for an anchor tweet that describes the mental state of the user. All the tweets published in a duration of one month from the anchor tweet were collected. They labeled a user as depressed if their anchor tweets satisfy the following regular expression “(I’m/I was/I am/I’ve been) diagnosed depression.” Thus, they obtained 1402 depressed users and 292564 tweets within one month. Similarly, the users who did not post any tweet having the word “depress,” were labeled as nondepressed.

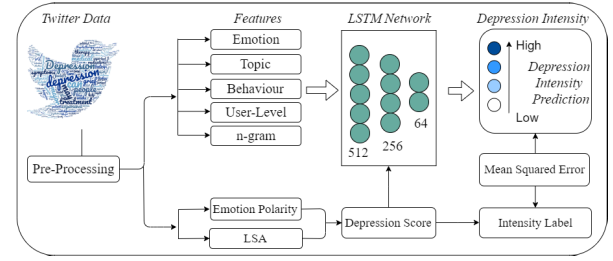


Fig. 3. Overall pipeline of the proposed method.

In this way, a nondepression dataset was also created. As the dataset is sparsely labeled into binary classes of depressed and nondepressed, these labels are not directly useful for depression intensity estimation. Therefore, we develop our own relabeling technique for the same dataset. Firstly, we compute a depression score based on the emotion polarity of tweets and LSA, as explained below. The scores are then mapped into four intensity categories.¹

1) *Emotion Polarity*: We use the NLTK library² to look into the compound polarity of tweets and assign an emotion polarity score in $[-1,1]$ scale to each user based on their tweets. NLTK mainly uses Valence Aware Dictionary and sEntiment Reasoner (VADER) algorithm output for sentiment scores having four classes of sentiments (i.e., neg: Negative, neu: Neutral, pos: Positive, and compound: Compound). Here, the compound is the aggregated score, which is a normalized score in the $[-1,1]$ range, approximated by the max expected value. As we label the depression intensity, we consider negative polarity as a positive value (higher overall score) and vice versa.

2) *LSA*: We extract all the keywords associated with the depression keyword in the latent semantic indexing (LSI) graph.³ It is an indexing method that uses the singular value decomposition (SVD) technique to identify patterns in a text document. LSI is mainly based on the principle that words used in a similar context tend to have similar semantic meanings. It can correlate semantically related terms. Table I shows the broad and narrow LSI keywords related to “depression” and Fig. 2 visualizes the closely related keywords. Words that are more generic abstract of the given keyword are called board terms and ones that are more specific are called narrow terms. Using the extracted keywords, we compute a *semantic score* for each user using the following (1), by taking the weighted sum of normalized broad term hit score $h(b)$ and normalized narrow term hit score $h(n)$, to make it range in $[0,1]$ scale. $h(b)$ and $h(n)$ are the normalized counts of hits of broad terms and narrow terms (related to depression in the LSI graph), respectively, with the user tweets. More weight can be given to the broad terms and less weights to the narrow terms

$$\text{Semantic score} = \alpha \times h(b) + (1 - \alpha) \times h(n). \quad (1)$$

¹This number can be changed depending upon application requirement. For our study, we keep it four.

²<http://www.nltk.org/howto/sentiment.html>

³<https://www.twinword.com/ideas/graph/depression/>

TABLE I
BROAD AND NARROW TERMS CORRESPONDING TO THE KEYWORD “DEPRESSION”

LSI Keywords	
Broad Terms	Narrow Terms
unhappiness, sadness, pushing, push, psychological state, psychological condition, mental state, mental condition, major affective disorder, incurvature, incurvation, formation, emotional disturbance, emotional disorder, economic condition, crisis, concavity, stress, anger, breakdown, cry, stress, affective disorder	wrinkle, depression, psychotic depression, prick, oppressiveness, oppression, oceanic abyss, neurotic depression, melancholy, melancholia, megrims, major depressive episode, lowland, low spirits, line, incision, hollow, hole, helplessness, heartsickness, groove, furrow, funk, fissure, exogenous depression, endogenous depression, dysthymic depression, dysthymia, dysphoria, droop, disconsolateness, dip, dimple, despondency, despondence, dent, demoralization, demoralisation, dejection, deep, crinkle, crevice, crease, crater, cranny, crack, click, chap,

3) *Depression Score*: We compute the final depression score by adding up the users’ emotion polarity and semantic scores together and taking their min-max normalization value, as shown in (2), where min and max are the minimum and maximum values of the summation of polarity and semantic scores. This way the user depression score lies in the range [0-1]

$$\text{Depression score} = \frac{(\text{polarity score} + \text{semantic score}) - \min}{\max - \min} \quad (2)$$

B. Data Preprocessing

We preprocess the data using the NLTK toolbox⁴ [46] and concatenate the user tweets for feature extraction. Our preprocessing step includes the following tasks. We remove the Emojis in Tweets’ texts as Emojis are not compatible with many text processing algorithms. All punctuation, articles, and special characters are removed. Before any processing, we tokenize the sentences. We used the Porter Stemmer [47] for stemming the tweets and WordNet Lemmatiser [48] to lemmatize them. We also process the irregular words generally mentioned on social media. These may be typographical mistakes or abbreviations of common words. Similar to [19], we leverage a word2vec model [49] trained on Twitter data for this task.

C. Feature Extraction

We aim to detect and analyze depressed users from their online behaviors. In the case of offline behaviors, there are prominent definitions in clinical depression criteria, which are widely used in depression diagnosis. Similarly, we extract some common online behaviors of social media users. From the literature, it is evident that emotion, event, online behavior, user-specific features, and use of depression-related texts, are relevant to depression intensity estimation. We consider the following features to describe each user.

1) *Emotional Features*: We consider 12-D emotion-related features, as defined in the following.

- 1) Negative words count of one dimension in tweets with the help of LIWC [42].
- 2) Emotion features both at the sentence level (4-D) and word level (4-D).⁵

3) *Valence arousal dominance (VAD) feature* [50] corresponding to the tweets (3-D). Basically, emotion is measured in 3-D plane called valance, arousal, and dominance. Valance represents the negative to positive effects of emotion. Arousal represents emotional intensity. Dominance represents control against emotional stimuli.

2) *Topical/Event Features*: We also consider events that trigger depression. For topic analysis, we use the latent Dirichlet allocation (LDA) technique [51]. We extract a 25-D topic feature in this regard.

3) *Online Behavioral Features*: The online behavior of users is captured by the following features.

- 1) Number of tweets (it is a 2-D vector which includes both user-level and within the time span value).
- 2) Social interactions (it is a 1-D vector which includes a number of retweets/comments per tweet. We summed up the value per tweet.).
- 3) Posting behaviors (it is a 2-D vector that includes a number of tweets per hour and the ratio of tweets per day to the total number of tweets).

4) *User Level Feature*: We use the following user attributes: “id,” “id str,” “name,” “screen name,” “location,” “profile location,” “description,” “url,” “entities,” “protected,” “followers count,” “friends count,” “listed count,” “created at,” “favourites count,” “utc offset,” “time zone,” “geo enabled,” “verified,” “statuses count,” “language,” “status,” “contributors enabled,” “is translator,” “is translation enabled,” “profile background color,” “profile background image url,” “profile background image url https,” “profile background tile,” “profile image url,” “profile image url https,” “profile banner url,” “profile link color,” “profile sidebar border color,” “profile sidebar fill color,” “profile text color,” “profile use background image,” “has extended profile,” “default profile,” “default profile image,” “following,” “follow request sent,” “notifications,” “translator type,” etc. Among these attributes, we do the following processing. The dates, url’s, and id’s are removed. The attributes having true/false values are converted to 0/1. For attributes having unique values, we obtain those unique sets. We give frequency-wise weightage to these attributes. For example, “time zone.” For text attributes like “status,” we extract the word2vec embedding corresponding to them. For the rest of the attributes having integer values, we normalize the column via min–max normalization. Finally, we get 334-D vector corresponding to the user-level features.

5) *Depression Related n-Gram*: We extract the *n*-grams related to depression keywords and their embedding as

⁴<https://www.datacamp.com/community/tutorials/stemming-lemmatization-python>

⁵<http://www.nltk.org/howto/sentiment.html>

features. We have used “depression,” “anxiety,” “stress,” “unhappy,” and “sad” keywords to extract the n -grams. For simplicity, we consider only unigrams, bigrams, and trigrams. We extract the following two features in this category.

- 1) We count the number of n -grams (1-D feature).
- 2) We extract 150-D word2vec embedding of these n -grams and used it as a feature.

6) *Overall Feature Set*: Thus, overall, we have 527 [12 (emotion) + 25 (topic level) + 5 (online behavior) + 334 (user level) + 151 (depression related n -gram)]-dimensional feature vectors corresponding to each user, as an input to the network.

D. Proposed Learning Model

We adopt LSTM networks [52], which are capable of learning long-term dependencies. We develop a three-layer architecture in the proposed method, as shown in Fig. 3. The details of the network architecture are shown in Table II and discussed later.

We use Swish [53] as an activation function instead of ReLU, defined as $f(x) = x \cdot \text{sigmoid}(x)$. Some of its properties (like unbounded above and bounded below) are similar to ReLU and some (like smooth and nonmonotonic) are different. Bounded below causes regularization and unbounded above speeds up the training process as it does not reach near-zero gradients. The main advantage of Swish lies in its self-gating property, i.e., it takes only one scalar value instead of multiple gating inputs. In short, it belongs in between ReLU and linear activation region. As a result of this, it possesses the positive points of both linear activation and ReLU. Considering a variation in Swish [54] that is $f(x) = 2x \cdot \text{sigmoid}(\beta \times x)$, where β is a learnable parameter. It can also be observed that, if $\beta = 0$ the remaining part becomes $1/2$ that is $f(x)$ becomes linear. Similarly, if the β is very high, the sigmoid part behaves like a binary activation (0 for $x < 0$ and 1 for $x > 0$). Therefore, the swish activation function ($\beta = 1$) incorporates a smooth transition between these two extremes.

V. EXPERIMENTS

We conduct extensive experiments and evaluate the performance of our method in several ways. First, we compare our results with the baseline intensity estimation models [support vector machine (SVM), deep neural network (DNN), and gated recurrent unit (GRU)] and then also compare with multimodal dictionary learning (MDL) in two different ways, as detailed in the following. In both cases, we outperform the others. We also evaluate the effectiveness of our relabeling technique and achieve good accuracy.

A. Methods for Comparison

SVM: We use an SVM regressor, which learns based on the separating hyperplane technique. Given labeled training data, the algorithm outputs an optimal hyperplane that can be able to capture data distribution.

DNN: We use a shallow DNN to predict the depression intensity as another baseline. The DNN layers are made

TABLE II
DNN, GRU AND LSTM ARCHITECTURE

DNN Layers	DNN Details	LSTM/GRU Layers	LSTM/GRU Details
Input	input layer	Input	input layer
Dense	512	LSTM/GRU	512
Activation	Relu/Swish	Activation	Relu/Swish
Dense	1024	LSTM/GRU	256
Activation	Relu/Swish	Activation	Relu/Swish
Dense	4096	LSTM/GRU	64
Activation	Relu/Swish	Activation	Relu/Swish
Depression (Sigmoid)	1	Depression (Sigmoid)	1

of hidden nodes. A hidden node is patterned on neurons corresponding to the human brain, which activates when it encounters sufficient stimuli. These nodes take input from the data with a set of weights which decides the significance of that given input. These input-weight products are summed and passed through an activation function.

GRU: The neural networks need to remember relevant patterns from the data. The major shortcoming of DNN is that it cannot remember any past states, especially for temporal data. Recurrent neural networks address this issue. RNNs have loops in them which allow any information to persist. GRU uses update gate and reset gate to remember/forget the previous states. These are two vectors that determine what information will be passed to the output.

MDL: MDL, proposed by Shen *et al.* [19], is based on features of six depression-related groups, covering both clinical depression criteria and online behaviors on social media. Their multimodal depressive dictionary learning model considers the problem as a binary classification task and identifies whether a user is depressed or not. In contrast, our work predicts the intensity of depression of a user, which makes a direct comparison difficult. Therefore, we compare with MDL in two different ways. First, MDL is adapted with our labeling technique for depression intensity estimation and compared with our proposed method and then our method is adapted for binary classification and compared with MDL.

B. Experimental Setup

We trained the learning methods with the following settings. For training an SVM, we use Gaussian kernel-based SVM. To train DNN, GRU, and LSTM, we use an SGD optimizer with cross-entropy as a loss function having a 0.01 learning rate with 0.9 momentum. The three-layered network details are given in Table II. All our experiments are performed on an intel core i7 processor having CPU 4.20-GHz frequency and 32-GB RAM. During training, we used Titan Xp GPU.

C. Evaluation Metrics

We evaluate the performance of the different models for depression intensity estimation in terms of mean squared error (MSE), which measures the average of the squares of errors. To evaluate the performance for binary classification, *accuracy* in terms of percentage is used as the evaluation metric. We perform fivefold cross-validation for testing and take the average over the fivefold.

TABLE III
RESULTS OF DEPRESSION INTENSITY ESTIMATION WITH SVM, DNN, GRU AND LSTM

Network Architecture	Feature	MSE (with relu)	MSE (with swish)
SVM	Emotion + Topic / Event + Online Behaviour + User level + n-gram	1.76	1.76
DNN	Emotion + Topic/Event + Online Behaviour + User level + n-gram	1.71	1.69
GRU	Emotion + Topic/Event + Online Behaviour + User level + n-gram	1.48	1.46
LSTM	Emotion	1.50	1.50
	Topic/Event	1.72	1.72
	Online Behaviour	1.68	1.66
	User level	1.80	1.80
	n-gram	1.58	1.56
	User level + Emotion	1.53	1.51
	User level + Topic/Event	1.67	1.67
	User level + Online Behaviour	1.62	1.61
	User level + n-gram	1.55	1.55
	Emotion + Topic / Event + Online Behaviour + User level + n-gram	1.43	1.42

TABLE IV
PERFORMANCE COMPARISON WITH MDL [19]: % MEASURES THE BINARY CLASSIFICATION ACCURACY AND MSE MEASURES THE INTENSITY PREDICTION ERROR

Binary classification method	Accuracy (%)	Intensity prediction method	MSE
MDL	84.80	MDL adapted with our labelling	1.88
Proposed method adapted with MDL labelling and fine-tuned	87.14	Proposed method	1.42

TABLE V
ACCURACY OF RELABELLING (TP: TRUE POSITIVES, TN: TRUE NEGATIVES)

Dataset	Users					
	[19]	Ours	TP	TN	Accuracy (severe)	Accuracy (moderate + severe)
Depressed	1,402	1,280	1,253	27	91.92%	100%
Non-depressed	5,160	540	NA	NA	NA	NA
		1,173				
		3,569				

D. Results of Depression Intensity Estimation

The quantitative results of our method are shown in Table III. Observe that the proposed model including the complete feature set trained on the LSTM network outperforms the other baseline models, achieving the lowest MSE of 1.42. By using the swish activation function the performance of the network increases slightly. We also trained our model with feature subsets for further analysis. In the case of feature-level prediction, the emotion feature is performing better as compared to others, as depression and emotion are interrelated. On the contrary, user-level information does not provide the important information corresponding to the depression intensity due to its diverse nature. As we have exploited user-level depression intensity, we combined each feature set with the user-level feature and run a similar LSTM network. Observed that the “user + emotion” feature performs better than other attributes. We also compare with MDL [19] to predict depression intensity (in [0-1] scale), by adapting MDL with our labeling technique. As shown in Table IV, clearly, our LSTM based proposed method, achieving an MSE of 1.42, outperforms adapted-MDL, achieving 1.88 MSE.

E. Representation Learning

In order to further evaluate whether our network learned a generalized representation of depression feature or not, we adapt our method for binary classification, as to whether a user is depressed or not and compare with [19]. We outperform their method by more than 2% of accuracy. The results are

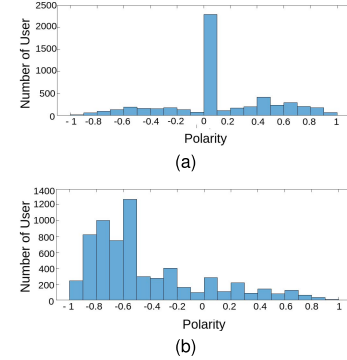


Fig. 4. Emotional polarity of users. (a) Nondepressed user. (b) Depressed user.

shown in Table IV. For a fair comparison, we use the binary labeling technique of [19], with our feature set and model to predict the depression label (i.e., whether the user is depressed or not). We extract the features from weights trained for depression intensity prediction and fine-tune for the binary classification. We fine-tuned our data on the Twitter depression data [19]. We added an FC layer having 128 nodes with ReLU activation to fine-tune the network. The learning rate was set to 0.001 with an SGD optimizer. For this task, we froze the pretrained weights of the network and fine-tuned the rest part.

F. Effectiveness of Our Relabeling Technique

The effectiveness of our labeling technique is illustrated in Table V. We divided the depression score [0-1] into four categories and evaluated these statistics by comparing them with the labels (depressed/nondepressed) already defined in the dataset. We achieve 91.92% accuracy for our severely depressed category and 100% accuracy by combining the *severely* and the *moderate-severely* depressed categories.

G. Online Behavior Analysis

Table VI shows some of the sample tweets from users identified at the different levels of depression. Similar to [19], we observe the following traits of depressed Twitter users.

TABLE VI
DEPRESSION INTENSITY RELATED TWEETS

Depression Intensities	Tweet Examples
no or minimal depression	I swear this person is fake trying to play a prank or something on me. I'm glad I'm not the only one who thinks this
mild depression	I hate how it's so easy for me to get into a bad mood i hateeee my friends Am I A Hypochondriac?: Hello. I have been diagnosed with mild depression and have an anxiety problem.
moderate depression	A friend says my life needs a restart. to delete Vet Medicine this time around go by way of the arts. I'm still lavving. Self Help Tips For Dealing W\depression & Mental Illness?
severe depression	I'm diagnosed somewhere between schizophrenia, bipolar disorder and severe psychotic depression. I need a doctor. Maybe some ropes too. Either to tie **** up or to tie a noose. I haven't decided yet. I deserve a golden ticket for my last tweet to ***.

Usage of Negative Words: The analysis of emotional polarity is analyzed in Fig. 4 in $[-1,1]$ scale where -1 means negative and $+1$ means positive. From the above-mentioned figure, it is observed that most of the nondepressed user does not use negative words in the post. On the other hand, depressed users mostly use negative words. Frequency wise top-3 negative words are “depression,” “stress,” and “sad.”

Posting Time: Depressed users (around 53%) mostly post tweets during the night. It indicates that tend to be insomniac. Lustberg and Reynolds III [55] also studied that depression symptoms tend to worsen during the night. According to their statistics, this can happen to eight out of ten people suffering from depression.

Linguistic Style: Similar to the work done by De Choudhury *et al.* [25], we analyzed articles, auxiliary verbs, conjunctions, adverbs, personal pronouns, prepositions, and negation. It is observed that depressed users use personal pronouns for indicating their suppressed nature.

Influence of Other Attributes: We also observed a few other attributes that trigger depression. Few personal events keywords (“work pressure,” “divorce,” “break-up,” etc.) are observed in their timeline. There may exist an event consequence relation between them. It is also observed that depressed users think it is a convenient platform to share feelings, get attention, or express their emotional state, especially feelings of helplessness.

VI. CONCLUSION

Depression is a pressing issue of our society today, affecting more than 264 million people globally and still increasing. It may become a serious health issue when it lasts more than two weeks with moderate or severe intensity. It can even lead to suicide if a depressed person is not receiving proper treatment. The situation is worst in countries like India, China, and the USA. During the ongoing COVID-19 pandemic and frequent lockdowns, mental health has become an important concern. In this article, we addressed the problem of early-stage depression detection from users' tweeting behavior. We proposed a deep learning method to estimate the intensity of depression by leveraging social media data. This work aims to make timely depression intensity estimation from social media, in order to aid in a proper treatment according to the depression level. We developed a relabeling technique for a benchmark depression dataset in a self-supervised manner,

designed a rich set of discriminative depression-related features for users, and proposed an LSTM network to detect depressed users of different levels on Twitter. We validated the performance of our method by conducting extensive experiments on a standard dataset and outperformed the other alternatives for intensity estimation. Our method adapted for binary classification outperforms the existing binary classification method by more than 2% of accuracy. This research leads to several promising future directions. It would be interesting to explore the social network structure and users locations for identifying the propagation of depression among the social communities. This research is also expected to provide further insights into depression research in the affective computing community. Moreover, as depression has become a pressing issue in the current pandemic scenario, we expect that the proposed research will lead to the development of automatic preliminary assessment methods based on social data.

REFERENCES

- [1] F. Hao, G. Pang, Y. Wu, Z. Pi, L. Xia, and G. Min, “Providing appropriate social support to prevention of depression for highly anxious sufferers,” *IEEE Trans. Comput. Social Syst.*, vol. 6, no. 5, pp. 879–887, Oct. 2019.
- [2] (2020). *Depression*. [Online]. Available: <https://www.who.int/news-room/fact-sheets/detail/depression>
- [3] J. A. Russell, “A circumplex model of affect,” *J. Personality Social Psychol.*, vol. 39, no. 6, p. 1161, Dec. 1980.
- [4] (2018). *India is the Most Depressed Country in the World*. [Online]. Available: <https://www.indiatoday.in/education-today/gk-current-affairs/story/india-is-the-most-depressed-country-in-the-world-mental-health-day-2018-1360096-2018-10-10>
- [5] R. I. Shader, “COVID-19 and depression,” *Clin. Therapeutics*, vol. 42, no. 6, pp. 962–963, 2020.
- [6] S. Pappa, V. Ntella, T. Giannakas, V. G. Giannakoulis, E. Papoutsis, and P. Katsaounou, “Prevalence of depression, anxiety, and insomnia among healthcare workers during the COVID-19 pandemic: A systematic review and meta-analysis,” *Brain, Behav., Immunity*, vol. 88, pp. 901–907, Aug. 2020.
- [7] J. Kanter and K. Manbeck. (2020). *Covid-19 Could Lead to an Epidemic of Clinical Depression, and the Health Care System isn't Ready for That, Either*. [Online]. Available: <https://theconversation.com/covid-19-could-lead-to-an-epidemic-of-clinical-depression-and-the-health-care-system-isnt-ready-for-that-either-134528>
- [8] D. Dresden. (2020). *Caring For Someone With Depression During The COVID-19 Pandemic*. [Online]. Available: <https://www.medicalnewstoday.com/articles/how-to-care-for-someone-with-depression>
- [9] C. Koons. (2020). *The Mental-Health-Care System isn't Ready for COVID-19 Either*. [Online]. Available: <https://www.bloomberg.com/news/articles/2020-04-01/the-u-s-mental-health-care-system-isnt-ready-for-coronavirus>
- [10] G. Coppersmith, M. Dredze, and C. Harman, “Quantifying mental health signals in Twitter,” in *Proc. Workshop Comput. Linguistics Clin. Psychol., Linguistic Signal to Clin. Reality*, 2014, pp. 51–60.

- [11] H. Lin, J. Jia, L. Nie, G. Shen, and T.-S. Chua, "What does social media say about your stress?" in *Proc. IJCAI*, 2016, pp. 3775–3781.
- [12] M. Akbari, X. Hu, N. Liqiang, and T.-S. Chua, "From tweets to wellness: Wellness event detection from Twitter streams," in *Proc. AAAI*, 2016, pp. 87–93.
- [13] T. Anwar and M. Abulaish, "Identifying cliques in dark Web forums—An agglomerative clustering approach," in *Proc. IEEE Int. Conf. Intell. Secur. Inform.*, Jun. 2012, pp. 171–173.
- [14] T. Anwar and M. Abulaish, "A social graph based text mining framework for chat log investigation," *Digit. Invest.*, vol. 11, no. 4, pp. 349–362, Dec. 2014.
- [15] T. Anwar and M. Abulaish, "Ranking radically influential Web forum users," *IEEE Trans. Inf. Forensics Security*, vol. 10, no. 6, pp. 1289–1298, Jun. 2015.
- [16] T. Anwar, K. Liao, A. Goyal, T. Sellis, A. S. M. Kayes, and H. Shen, "Inferring location types with geo-social-temporal pattern mining," *IEEE Access*, vol. 8, pp. 154789–154799, 2020.
- [17] M. Park, C. Cha, and M. Cha, "Depressive moods of users portrayed in twitter," in *Proc. ACM SIGKDD Workshop Healthcare Informat. (HI-KDD)*, 2012, pp. 1–8.
- [18] A. T. Beck, "An inventory for measuring depression," *Arch. Gen. Psychiatry*, vol. 4, no. 6, p. 561, Jun. 1961.
- [19] G. Shen *et al.*, "Depression detection via harvesting social media: A multimodal dictionary learning solution," in *Proc. 27th Int. Joint Conf. Artif. Intell.*, Aug. 2017, pp. 3838–3844.
- [20] T. Shen *et al.*, "Cross-domain depression detection via harvesting social media," in *Proc. 27th Int. Joint Conf. Artif. Intell.*, Jul. 2018, pp. 1611–1617.
- [21] A. T. Beck, R. A. Steer, R. Ball, and W. F. Ranieri, "Comparison of beck depression inventories-IA and-II in psychiatric outpatients," *J. Personality Assessment*, vol. 67, no. 3, pp. 588–597, Dec. 1996.
- [22] L. S. Radloff, "The CES-D scale: A self-report depression scale for research in the general population," *Appl. Psychol. Meas.*, vol. 1, no. 3, pp. 385–401, Jun. 1977.
- [23] O. Whooley, "Diagnostic and statistical manual of mental disorders (DSM)," in *The Wiley Blackwell Encyclopedia of Health, Illness, Behavior, and Society*. 2014.
- [24] M. Park, D. W. McDonald, and M. Cha, "Perception differences between the depressed and non-depressed users in Twitter," in *Proc. AAAI Conf. Weblogs Social Media*, 2013, pp. 476–485.
- [25] M. De Choudhury, M. Gamon, S. Counts, and E. Horvitz, "Predicting depression via social media," in *Proc. 7th Int. AAAI Conf. Weblogs Social Media*, 2013, pp. 128–137.
- [26] P. Resnik, W. Armstrong, L. Claudino, T. Nguyen, V.-A. Nguyen, and J. Boyd-Graber, "Beyond LDA: Exploring supervised topic modeling for depression-related language in Twitter," in *Proc. 2nd Workshop Comput. Linguistics Clin. Psychol., Linguistic Signal Clin. Reality*, 2015, pp. 99–107.
- [27] S. Tsugawa, Y. Kikuchi, F. Kishino, K. Nakajima, Y. Itoh, and H. Ohsaki, "Recognizing depression from Twitter activity," in *Proc. 33rd Annu. ACM Conf. Hum. Factors Comput. Syst.*, Apr. 2015, pp. 3187–3196.
- [28] Y. Xue, Q. Li, L. Feng, G. D. Clifford, and D. A. Clifton, "Towards a micro-blog platform for sensing and easing adolescent psychological pressures," in *Proc. ACM Conf. Pervas. Ubiquitous Comput. Adjunct Publication*, Sep. 2013, pp. 215–218.
- [29] Y. Xue, Q. Li, L. Jin, L. Feng, D. A. Clifton, and G. D. Clifford, "Detecting adolescent psychological pressures from micro-blog," in *Proc. Int. Conf. Health Inf. Sci.* Melbourne, VIC, Australia: Springer, 2014, pp. 83–94.
- [30] R. Xu and Q. Zhang, "Understanding online health groups for depression: Social network and linguistic perspectives," *J. Med. Internet Res.*, vol. 18, no. 3, p. e63, Mar. 2016.
- [31] F. Sadeque, D. Xu, and S. Bethard, "Measuring the latency of depression detection in social media," in *Proc. 11th ACM Int. Conf. Web Search Data Mining*, Feb. 2018, pp. 495–503.
- [32] M. Trotszek, S. Koitka, and C. M. Friedrich, "Utilizing neural networks and linguistic metadata for early detection of depression indications in text sequences," *IEEE Trans. Knowl. Data Eng.*, vol. 32, no. 3, pp. 588–601, Mar. 2020.
- [33] T. Cai, J. Li, A. S. Mian, R. Li, T. Sellis, and J. X. Yu, "Target-aware holistic influence maximization in spatial social networks," *IEEE Trans. Knowl. Data Eng.*, early access, Jun. 17, 2020, doi: [10.1109/TKDE.2020.3003047](https://doi.org/10.1109/TKDE.2020.3003047).
- [34] J. Li, X. Wang, K. Deng, X. Yang, T. Sellis, and J. X. Yu, "Most influential community search over large social networks," in *Proc. IEEE 33rd Int. Conf. Data Eng. (ICDE)*, Apr. 2017, pp. 871–882.
- [35] J. Li, C. Liu, J. X. Yu, Y. Chen, T. Sellis, and J. S. Culpepper, "Personalized influential topic search via social network summarization," *IEEE Trans. Knowl. Data Eng.*, vol. 28, no. 7, pp. 1820–1834, Jul. 2016.
- [36] J. Li and C. Cardie, "Timeline generation: Tracking individuals on Twitter," in *Proc. 23rd Int. Conf. World Wide Web (WWW)*, 2014, pp. 643–652.
- [37] J. Li, A. Ritter, C. Cardie, and E. Hovy, "Major life event extraction from Twitter based on congratulations/condolences speech acts," in *Proc. Conf. Empirical Methods Natural Lang. Process. (EMNLP)*, 2014, pp. 1997–2007.
- [38] L. S. Greenberg and J. C. Watson, *Emotion-Focused Therapy for Depression*. Washington, DC, USA: American Psychological Association, 2006.
- [39] A. Pak and P. Paroubek, "Twitter for sentiment analysis: When language resources are not available," in *Proc. 22nd Int. Workshop Database Expert Syst. Appl.*, Aug. 2011, pp. 111–115.
- [40] C. Strapparava and A. Valitutti, "Wordnet affect: An affective extension of WordNet," in *Proc. LREC*, vol. 4, 2004, p. 40.
- [41] S. Ghosh, M. Chollet, E. Laksana, L.-P. Morency, and S. Scherer, "Affect-LM: A neural language model for customizable affective text generation," in *Proc. 55th Annu. Meeting Assoc. Comput. Linguistics*, vol. 1, 2017, pp. 634–642.
- [42] J. W. Pennebaker, M. E. Francis, and R. J. Booth, "Linguistic inquiry and word count: LIWC 2001," *Mahway, Lawrence Erlbaum Associates*, vol. 71, no. 2001, p. 2001. 2001.
- [43] X. Rong, "Word2vec parameter learning explained," 2014, *arXiv:1411.2738*. [Online]. Available: <http://arxiv.org/abs/1411.2738>
- [44] N. Asghar, P. Poupard, J. Hoey, X. Jiang, and L. Mou, "Affective neural response generation," in *Proc. Eur. Conf. Inf. Retr.* Grenoble, France: Springer, 2018, pp. 154–166.
- [45] *What is Depression?* Accessed: May 12, 2020. [Online]. Available: <https://www.psychiatry.org/patients-families/depression/what-is-depression>
- [46] E. Loper and S. Bird, "NLTK: The natural language toolkit," 2002, *arXiv:cs/0205028*. [Online]. Available: <https://arxiv.org/abs/cs/0205028>
- [47] M. F. Porter. (2001). *Snowball: A Language for Stemming Algorithms*. Accessed: May 12, 2020. [Online]. Available: <http://snowball.tartarus.org/texts/introduction.html>
- [48] S. Poria, A. Gelbukh, E. Cambria, P. Yang, A. Hussain, and T. Durrani, "Merging SenticNet and WordNet-affect emotion lists for sentiment analysis," in *Proc. IEEE 11th Int. Conf. Signal Process.*, Oct. 2012, pp. 1251–1255.
- [49] T. Baldwin, M.-C. de Marneffe, B. Han, Y.-B. Kim, A. Ritter, and W. Xu, "Shared tasks of the 2015 workshop on noisy user-generated text: Twitter lexical normalization and named entity recognition," in *Proc. Workshop Noisy User-Generated Text*, 2015, pp. 126–135.
- [50] M. M. Bradley and P. J. Lang, "Affective norms for English words (ANEW): Instruction manual and affective ratings," Citeseer, Gainesville, FL, USA, Tech. Rep. C-1, 1999.
- [51] D. M. Blei, A. Y. Ng, and M. I. Jordan, "Latent Dirichlet allocation," *J. Mach. Learn. Res.*, vol. 3, pp. 993–1022, Mar. 2003.
- [52] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Comput.*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [53] P. Ramachandran, B. Zoph, and Q. V. Le, "Searching for activation functions," 2017, *arXiv:1710.05941*. [Online]. Available: <http://arxiv.org/abs/1710.05941>
- [54] E. Alcaide, "E-swish: Adjusting activations to different network depths," 2018, *arXiv:1801.07145*. [Online]. Available: <http://arxiv.org/abs/1801.07145>
- [55] L. Lustberg and C. F. Reynolds, "Depression and insomnia: Questions of cause and effect," *Sleep Med. Rev.*, vol. 4, no. 3, pp. 253–262, Jun. 2000.