

Received March 11, 2019, accepted March 21, 2019, date of publication April 4, 2019, date of current version April 16, 2019.

Digital Object Identifier 10.1109/ACCESS.2019.2909180

Detection of Depression-Related Posts in Reddit Social Media Forum

MICHAEL M. TADESSE¹, HONGFEI LIN¹, BO XU¹, AND LIANG YANG

Faculty of Electronic Information and Electrical Engineering, Dalian University of Technology, Dalian 116024, China

Corresponding author: Hongfei Lin (hflin@dlut.edu.cn)

This work was supported in part by the Natural Science Foundation of China under Grant 61632011, Grant 61572102, Grant 61602078, Grant 61572098, in part by the Fundamental Research Funds for the Central Universities under Grant DUT18ZD102, in part by the Ministry of Education Humanities and Social Science Project under Grant 19YJCZH199, in part by the China Postdoctoral Science Foundation under Grant 2018M641691, and in part by the National Key Research Development Program of China under Grant 2016YFB1001103.

ABSTRACT Depression is viewed as the largest contributor to global disability and a major reason for suicide. It has an impact on the language usage reflected in the written text. The key objective of our study is to examine Reddit users' posts to detect any factors that may reveal the depression attitudes of relevant online users. For such purpose, we employ the Natural Language Processing (NLP) techniques and machine learning approaches to train the data and evaluate the efficiency of our proposed method. We identify a lexicon of terms that are more common among depressed accounts. The results show that our proposed method can significantly improve performance accuracy. The best single feature is bigram with the Support Vector Machine (SVM) classifier to detect depression with 80% accuracy and 0.80 F1 scores. The strength and effectiveness of the combined features (LIWC+LDA+bigram) are most successfully demonstrated with the Multilayer Perceptron (MLP) classifier resulting in the top performance for depression detection reaching 91% accuracy and 0.93 F1 scores. According to our study, better performance improvement can be achieved by proper feature selections and their multiple feature combinations.

INDEX TERMS Natural language processing, machine learning, Reddit, social networks, depression.

I. INTRODUCTION

Depression as a common mental health disorder has long been defined as a single disease with a set of diagnostic criteria. It often co-occurs with anxiety or other psychological and physical disorders; and has an impact on feelings and behaviour of the affected individuals [1]. According to the WHO study, there are 322 million people estimated to suffer from depression, equivalent to 4.4% of the global population. Nearly half of the in-risk individuals live in the South-East Asia (27%) and Western Pacific region (27%) including China and India. In many countries depression is still under-diagnosed and left without any adequate treatment which can lead into a serious self-perception and at its worst, to suicide [2]. In addition, the social stigma surrounding depression prevents many affected individuals from seeking an appropriate professional assistance. As a result, they turn to less formal resources such as social media.

The associate editor coordinating the review of this manuscript and approving it for publication was Vijay Mago.

With the development of Internet usage, people have started to share their experiences and challenges with mental health disorders through online forums, micro-blogs or tweets. Their online activities inspired many researchers to introduce new forms of potential health care solutions and methods for early depression detection systems. Using different Natural Language Processing (NLP) techniques and text classification approaches, they tried to succeed in a higher performance improvement. Some studies use single set features, such as bag of words (BOW) [3], [4], N-grams [5], LIWC [6] or LDA [7], [8] to identify depression in their posts. Some other papers compare the performance of individual features with various machine learning classifiers [9]–[12]. Recent studies examine the power of single features and their combinations such as N-grams+LIWC [13] or BOW+LDA and TF-IDF+LDA [14] to improve the accuracy results. They experiment with a smarter text pre-processing, and introduce different substitute words depending on the nature of the original string. For instance, Tyshchenko et al. [14] suggested categorizing the stop words and adding LIWC-like

word categories as an extra feature to an already designed method (BOW+TFIDF+LIWC). In addition, he applied multiple feature combinations to increase the performance using Convolutional Neural Networks (CNN) which consist of neurons with learnable weights and differ in terms of their layers. CNNs are very similar to simple feed-forward neural networks and state of the art method in the text and sentence classification tasks. A meta-analysis by Guntuku et al. [15] summarizes several iterations of depression detection tasks in computational linguistics. Another interesting review for mental health support and intervention in social media is written by Calvo et al. [16] who reviewed the taxonomy of data sources, NLP techniques and computational methods to detect various mental health applications.

Even with this significant progress, challenges still remain. This paper aims to search for a solution to a performance increase through a proper features selection and their multiple feature combinations. First, we choose the most beneficial linguistic features applied for depression identification to characterize the content of the posts. Second, we analyze the correlation significance, hidden topics and word frequency extracted from the text. Regarding the correlation, we focus on the LIWC dictionary and its three feature types (linguistic dimensions, psychological processes and personal concerns). For the topic examination, we choose the LDA method as one of the successful features. For the word frequency, we use unigrams and bigrams by leveraging the vectors based on TF-IDF scheme. Finally, we set five text classifying techniques and conduct their execution using the extracted data to detect depression. We compare the performance results based on three single feature sets and their multiple feature combinations. In our experiment, we use data collected from the Reddit social media platform. It was chosen as the data source as it allows longer posts. Targeting technical approaches towards detection tasks, our paper follows the lines of Calvo et al. research [17].

Our study has four specific contributions: first, to examine the relationship between depression and user's language usage; second, to design three LIWC features for our specific research problem; third, to evaluate the power of N-grams probabilities, LIWC and LDA as single features for performance accuracy; fourth, to show the predictive power of both single and combined features with proposed classification approaches to achieve a higher performance in depression identification tasks.

The rest of the paper is organized as follows. In section II, we discuss related work in depression detection. In section III, we define the properties of the Reddit dataset. In section IV, we introduce the methodology and conduct data pre-processing followed by feature extraction. In section V, we compare and analyze the feature sets and examine the results as well as the most powerful machine learning technique for depression detection. We conclude our study and provide a direction for future work in section VI.

II. RELATED WORK

There are various types of studies examining the relationship between mental health and language usage to provide new insight into depression detection. Dating back to the earliest years of psychology, Sigmund Freud [18] wrote about Freudian slips or linguistic mistakes to reveal the secret thoughts and feelings of the writers. With the development of sociology and psycholinguistic theories, various approaches towards the relationship between depression and its language have been defined. For instance, according to Aaron Beck et al. [19]'s cognitive theory of depression, affected individuals tend to perceive themselves and their environment in mostly negative terms. They often express themselves through negatively valenced words and first-person pronouns. Their typical feature is self-preoccupation defined by Pyszvzinsky and Greenberg [20] which can develop into an extreme self-criticism stage. According to Durkheim's [21] social integration model, people suffering from depression often feel detached from their social life and have a difficulty to integrate into society.

These theories have motivated other researchers to come up with empirical support for their validity. For instance, Stirman and Pennebaker [25] compared the word usage of 300 poems written by 9 suicidal and 9 non-suicidal writers in three different periods of their lives. The results show that the suicidal poets used more first-person singular pronouns (I, me or we). Similar experiment was done by Rude et al. [26] who examined the linguistic patterns of the essays written by currently-depressed, formerly-depressed and never depressed college students. According to his results, depressed students used more negatively valenced words and less positive emotion words. Zinken et al. [27] studied the psychological relevance of syntactic structures to predict the improvement of depressive symptoms. He supposed that a written text might barely differ in its word usage; however, may differ in its syntactic structure, especially in the construction of relationships between the events. Analyzing a causation and insight words tasks, he found out that in the text written by depressed individuals there was a decreased use of complex syntax in comparison to non-depressed ones.

With the development of social media and Internet age, studies about depression and other mental health disorders have brought new challenges. Online domains such as Facebook, Twitter or Reddit have created a new platform for innovative research with a rich source of text data and social metadata to capture the users' behavioral tendencies. NLP techniques and various classifying approaches have been applied to evaluate the textual data and examine the impact of social networks on the mental health of the users. The data have been evaluated from different perspectives such as the text level, author level, and community level.

Twitter is one of the most popular social networking sites with almost 326 million active users and 90 million tweets publicly broadcasted to a large audience [28]. Many researchers have successfully utilized Twitter data

as a source of insights into the epidemiology of emotions, depression and other mental disorders of the users who tweet. De Choudhury et al. [29] used linguistic features to train a classifier to examine Twitter posts that indicated depression. Coppersmith et al. [6] looked for tweets that explicitly stated “I was just diagnosed with depression” sentences. Preotiuc-Pietro et al. [9] applied broader textual features such as LIWC, LDA and frequent 1-3 grams on the Twitter data to examine the personality of the users with self-declared post-traumatic stress (PTSD) disorders. His results show that the users suffering from PTSD were both older and more conscientious in comparison to depressed individuals. Since the language predictive of depression and PTSD had a large overlap with the language predictive of personality, the authors conclude that the users with a particular personality or demographic profiles tend to share their mental health diagnosis on social media, and thus the results may not generalize to other sources of autobiographical text. Resnik et al. [8] proved that the LDA model can uncover a meaningful and potentially useful latent structure for the automatic identification of important topics for depression detection. Tsugawa et al. [12] predicted depression from Twitter data in a Japanese sample where he showed that the features based on a topic modeling are useful in the tasks for recognizing depressive and suicidal users. Bentoni et al. [5] demonstrated the effectiveness of multi-task learning (MTL) models on mental health disorders with a limited amount of target data. He used feed-forward multi-layer perceptrons and feed-forward multi-task models trained to predict each task separately as well as to predict a set of conditions simultaneously. They experimented with a feed-forward network against independent logistic regression models to test if MTL would have performed well in the domain. Reece et al. [22] found out that the first stage of depression may be detectable from Twitter data several months prior to its diagnosis with 0.87 AUC of performance probability.

Facebook is another micro-blogging site with more than 2.2 billion users [28] providing a rich source of data for further research. For instance, Moreno et al. [30] evaluated the status updates of 200 Facebook users by using the references to “self-declared” diagnoses identified through “I feel hopeless” statements. Schwartz et al. [11] believed that depression severity is continuously distributed rather than dichotomous. In his experiment, he observed seasonal fluctuations of depression and found out that depression can be viewed as a continuous construct with constant changes reaching its highest degree during the winter months. Eichstaedt, J.C. et al. [24] study proved that it is possible to predict depression from online electronic medical records within the period up to six months achieving 0.72 AUC performance. He identified depressed patients from the records accessed through the Facebook statuses posted by the patients visiting the same emergency department. His findings show that depression is marked by linguistic predictors such as increased perceptual processes, references to sadness and discrepancies or greater negative emotions.

Reddit social media is widely used as an online discussion forum conducted through different communities or “subreddits”. Since it allows full anonymity of the users, it is often used for discussions about stigmatic topics. Choudhury et al. [31] analyzed the posts of Reddit users who wrote about mental health discourse and later shifted to discuss topics about suicidal ideation. This shift was predicted by the features such as self-concern, poor linguistic style, reduced social engagement, and expression of hopelessness or anxiety. Bagroy et al. [32] examined the potential of social media for studying the mental health of college students from over 100 universities. Based on the results, the amount of the posts connected with depression increased over the course of the academic year, especially at the universities with quarter-based schedules. Maupom et al. [7] implemented a system based on the topic extraction algorithm and simple neural networks extracting 30 latent topics in an unsupervised manner. His result shows that the limited number of users greatly hindered the predicting power of the MLP.

In recent past, shared tasks potentially applicable to different situations have become significantly popular in a wide research community. CLEF eRISK or the Conference and Labs of Evaluation Forum for Early Risk Prediction, is a public competition that allows researchers from multiple disciplines to join and collaborate on the creation of reusable benchmarks for evaluating early risk detection technologies employed in different areas such as health and safety [33]. The CLEF eRisk 2018 focused on early detection of signs of depression or anorexia. The corpora was built by Reddit textual data. Different machine learning techniques and feature engineering schemes were applied on the dataset to achieve the best results [4], [34], [35].

Non-English social media are quite unexplored areas in computational linguistics with their own conception of mental health and social stigma worthy of further study. For instance, Masuda et al. [36] examined data from Japanese Mixi social media platform to analyze the relationship between a user’s social network and suicide ideation. His findings show that the impact of the age, gender and number of friends on suicide ideation was small. Li et al. [37] analyzed the social attitude of Chinese Weibo users commenting on the posts of others who publicly broadcasted their intention to commit suicide. He found out that social stigma was widespread and words such as stupid, pathetic or deceitful were often used to make a comment on the posts of the in-risk individuals. Ramirez-Esparza et al. [38] showed that depressed users who wrote in Spanish were more likely to mention relational concerns rather than depressed one who wrote in English.

Apart from above mentioned social media, there are other online variables offering data sources for depression detection research. For instance, Nguyen et al. [10] experimented with a wide range of NLP techniques in LiveJournal social networking service to transform the text into a high-dimensional space and captured its topics and moods posted on online blogs. He aimed to present a good predictive

TABLE 1. Overview of recent depression detection studies in social media.

Data Source	Features	Methods	Best Results	References
Reddit	BOW, UMLS	Ada Boost & SVM	F1-score 0.75 & 0.98	Sayanta Paul et al. [4]
		MLP	F1-score 0.64	
	LIWC+N-grams	SVM	Acc. 0.82	JT Wolohan et al. [13]
Tweets	LIWC, sentiment, time series	RF	AUC 0.87	Reece et al. [22]
		LR	AUC 0.85	Preotiuc-Pietro et al. [9]
	LIWC, n-grams, topic modeling, sentiment	SVM	Acc. 0.69	Tsugawa et al. [23]
	N-grams, topic modeling	Neural network	AUC 0.76	Benton et al. [5]
Facebook	LIWC, N-grams, topic modeling	LR	AUC 0.72	Eichstaedt, J.C. et al. [24]
LiveJournal	LIWC, topics modeling, mood tags	Regression Models	Acc. 0.93	Nguyen et al. [10]
Blog posts	TFIDF, topic modeling, BOW	CNN	Acc. 0.78	Tyshchenko [14]

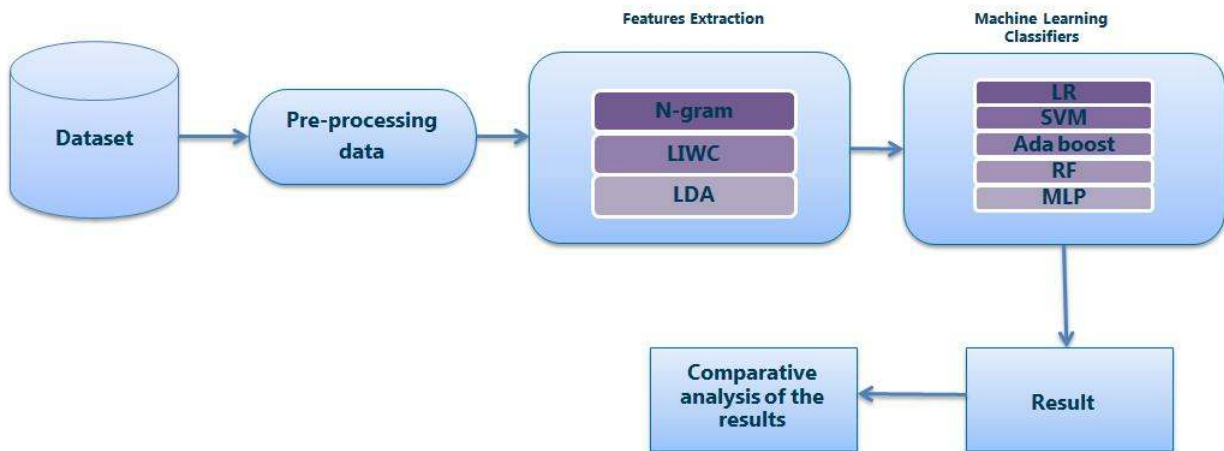


FIGURE 1. Depression Detection Framework.

validity in the depression classification between the clinical and control groups.

Table 1 provides an overview of research papers that have mined textual data for insights into the relationship between depression and language usage. The first column describes the source of data used in social media, Livejournal and blog posts. The features are based on a linguistic analysis of the text followed by different machine learning approaches to obtain the best performance for the particular experiment conducted by the authors mentioned in the last column.

III. EVALUATION DATASETS

To identify depression, we train our models on the dataset of Reddit users. The dataset was built by Inna Pirina et al. [39] and consists of a list of depressed and non-depressed users.

Since the users are often active in different subreddits, each group is created by a corresponding random number of messages covering different topics. The data corpus contains depression-indicative posts (1293) and standard posts (548). Depression-indicative posts are collected from relatively large subreddits devoted to depression, where depressed users seek support from an online community. Standard posts written by non-depressed users are collected from subreddits related to a family or friends. Table 2 demonstrates the words for the posts in both categories which are topically specific.

IV. METHODOLOGY

There are a growing number of methodologies to detect depression from the posts. In our study, we incorporate a

TABLE 2. Words frequently used in depression-indicative posts and standard posts.

Depression-Indicative Posts	Standard Posts
alone, break, blame, depressed, deserve better, deserve unhappy, die, escape, distraction, nobody, feel alone, feel depressed, felt pain, fuck don't, hate, hurt, loneliness, mine, myself, reject love, safe, shit, sucks, no job, painful, pressure, too worried, unsuccessful, ugly, uncomfortable, winter, worry, worth, wrong life	awesome, aunts, believe, beautiful, close, advice, cooking, cousins, don't care, encourage, family, logical person, got married, I do, better, mom, peace, parents, spend time, new friends, right, funny, need, thankfully, uncles, soul-friends, work, weekend, movie, potential, texted me, too good

technical description of approaches applied for depression identification using the NLP and text classifying techniques. The framework in Fig.1 consists of data-pre-processing, feature extraction followed by the machine learning classifiers, features analysis and experimental results.

A. DATA PRE-PROCESSING

We use the NLP tools to pre-process the dataset before it is proceeded to the feature selection and training stage. First, we use tokenization to divide the posts into individual tokens. Next, we remove all the URLs, punctuations and stop words which could lead into erratic results if stay ignored. Then we apply stemming in order to reduce the words to their root form and group similar words together.

B. FEATURES EXTRACTION

After data pre-processing, we feed our models with the features that reflect users' language habits in Reddit forums. To explore the users' linguistic usage in the posts, we employ the LIWC dictionary, LDA topics, and N-gram features. These text encoding methods are applied to encode the words to be proceeded by different classifiers.

N-gram modeling is used to examine the features from the posts. It is widely used in text mining and NLP as a feature for depression detection [9], [40] to calculate the probability of co-occurrence of each input sentence as a unigram and bigram. For n-gram modelling we use the Term frequency-inverse document frequency (TF-IDF) as a numeric statistic where the importance of a word with respect to each document in corpora is highlighted. The main goal of its usage is to scale down the impact of empirically less informative tokens, which occur frequently to give a space for the more informative words occurring in a small fraction. The word is ranked with greater TF-IDF value if it is present in a particular post and absent in other post [41]. In our study, we use TF-IDF vectorizer from the scikit-learn Python library [42] to extract 194,613 unigrams and bigrams. We remove all the stop words from the dataset and restrict the term-document matrix to 3000 most frequent unigrams and bigrams. In addition, we used Pointwise mutual information (PMI) [43] to filter infrequent bigrams.

LIWC, or the Linguistic Inquiry and Word Count dictionary, is widely used in computational linguistics as a source of

TABLE 3. Different types of approaches to text encoding methods.

Feature Type	Methods	Number of Selected Features
N-grams	unigram	3000
	bigram	2736
Linguistic dimensions Psychological Processes Personal Concern Processes	LIWC	68
Topic Modeling	LDA	70

features for psychological and psycholinguistic analysis [44]. It works as a baseline measure with a set of words and a behavioral link. It is often presented in several mental health projects [6], [11], [12], [45]–[48]. Table 3 describes different types of approaches to text encoding methods.

To accomplish our experiment, we extract 68 among 95 different features in view of psycholinguistic measures and change every depressive and non-depressive post into numerical values. This way we obtain the scores for three higher-level categories considering standard linguistic dimensions, psychological processes and personal concerns. The standard linguistic processes are one of the largest parts of the LIWC psycholinguistic vocabulary package. It was intended to quantify the words' usage in mentally significant classifications as well as for recognizing the connection between individuals in social co-operation. In our study, we first choose 9 linguistic features (articles, auxiliary verbs, adverbs, conjunctions, impersonal and personal pronouns, negations, prepositions and verbs) to characterize the users' text. Then we divide the Psychological processes into subcategories from which we used effective processes (anxiety, sadness, positive or negative emotion), biological processes (sexual, body, ingestion and health), social processes (family, friend, male, female), cognitive processes (cause, always, never, because), personal concerns (job, cook, cash, bury, kill), and time orientations (present, past, season). To examine the users' linguistic usage, we implement LIWC2015 dictionary [44] as the pre-defined category to measure all the textual content submitted by the users to extract lexico-syntactic features. We evaluate the correlation using the Pearson correlation coefficient r and also Benjamini-Hochberg selection method used in [24].

Topic modelling is an effective tool in computational linguistics to reduce the input of textual data feature space to a fixed number of topics [49]. Through the unsupervised text mining approach, hidden topics such as topics connected with anxiety and depression can be extracted from the selected documents. In comparison to LIWC, it is not created by a



FIGURE 2. Most Frequent N-grams in Depression-indicative Posts.

fixed set of pre-established words. However, it automatically generates the group of non-labelled words. The choice of words is based on a probability. As a result, each generated document deals with different topics that keep some link among each other. In our study, we examine the content of each post semantically connected with depression discussion session in Reddit. To derive the topic distributions for each post in the dataset, we used Latent Dirichlet Allocation (LDA) module. It is a probabilistic generative model for discretization of data collections helpful in discovering its underlying topic structures [50].

Based on our results, LDA model works best on the validation set when it is limited to 70 topics. For the topic selection we consider only the words that appear at least in more than 10 posts. We include every post as a single document that must be further tokenized and stemmed. This way allows us to compute the topics over the collection of documents to annotate them according to detected topics. Before we start the topic modelling process, all the stop words are removed. LDA implementation is provided by the Mallet toolkit [51].

C. TEXT CLASSIFICATION TECHNIQUES

To estimate the presence of depression, we employ classifying approaches to estimate the likelihood of depression within the users. The proposed framework is developed by using Logistic Regression, Support Vector Machine, Random Forest, Adaptive Boosting and Multilayer Perceptron classifier.

Logistic Regression (LR) is a linear classification approach used to estimate the probability occurrence of binary response based on one or more predictors [52], [53] and features.

Support Vector Machine (SVM) model is a representation of the examples as points in a highly dimensional space utilized for classification, where the points of the separate categories are widely divided. New examples are then mapped into the same space and predicted to belong to a category based on which side of the gap they fall [54].

Random Forest (RF) is an ensemble of decision tree classifiers trained with the bagging method where a combination of learning models increases the overall result [55].

Adaptive Boosting (AdaBoost) is an ensemble technique that can combine many weak classifiers into one strong

classifier [56]. It is widely used for binary class classification problems [57].

Multilayer Perceptron (MLP) is a special case of the artificial neural network often used for modelling complex relationships between the input and output layers [58]. Due to its multiple layers and non-linear activation it can distinguish the data that is not only non-linearly separable [59]. In our study, we applied the MLP method and two hidden layers with 4 and 16 perceptrons to fix for all the features in order to ensure a comparison consistency.

V. EXPERIMENTAL RESULTS

Before we analyze the classification results, we discuss several quantitative results. We define the words frequency, examine the correlation significance between the words and the features, and analyze the predictive power of LDA features.

A. FEATURES ANALYSIS

1) FREQUENCY AND PREDICTIVE POWER OF N-GRAM FEATURES

To compare the differences in the lexicon, we classify the entire labelled corpus of Reddit posts. To investigate the presence of depression, we compute the frequencies of all the unigrams and bigrams in both depression-indicative posts and standard posts. We select top 100 unigrams and bigrams for each category. Fig.2 and Fig.3 presents the top 100 unigrams and top 100 bigrams for each type of the posts generated in our study. The emerged words indicated by a high frequency are illustrated in both figures.

According to our results, the language predictors of depression in depression-indicative posts contain the words related to preoccupation with the self (I'm, I'm not), feelings (feel, feel like, make feel), greater discrepancies (want, wish, could), negative emotions (sad, wrong life, depression, I'm depressed, miserable), suicidal thoughts (want die, stop-stop, don't want, kill, mental illness), words of anger and hostility (shit, f***, hate), words of negation (no, not, no one, doe anyone, I'm not), interpersonal processes (lonely, feel alone), signs of hopelessness (pointless, anyone else, need help, end),



FIGURE 3. Most Frequent N-grams in Standard Posts.

signs of meaninglessness (pointless, empty, senseless) and present tense events (I'm going). As depression often affects psychomotor functions [60], we can find the words (tired, I'm tired or sleep) which reflect the symptoms of low energy, fatigue or inversely insomnia and hyperactivity. It is often expressed somatically through the bodily symptoms (my head, pain, hurt). In contrary to depression-indicative posts, unigrams and bigrams in standard posts contain the words describing the events happening rather in the past (time, month ago, year ago, last year), social relations (friend, best friend, family, friendship, mom) and advice seeking words (need advice, please help).

2) PREDICTIVE POWER OF LIWC FEATURES

We selected 68 out of 95 features to analyze the correlation between the textual data and the features themselves. We converted every depressive and non-depressive post into numerical values in view of psycholinguistic features resulting in correlation presented in the features extraction. The highest correlation is found in the Psychological processes (0.19) followed by the Linguistic dimensions (0.17) and Personal concerns (0.16). In subcategories of the Psychological processes, we examined the Affective processes (kill, worthless, cry) with the highest correlation results (0.19), especially for negative (0.19) and positive emotions (0.09). It is followed by the Social processes (brother, friend, mom) (0.17) with the results similar to the Linguistic dimensions (0.17). Based on the Personal concerns, the users like to talk about topics such as work (0.16), money (0.14) or death (0.11).

Concerning the mental focus of depressed and non-depressed users, the results show that depressed individuals use more self-oriented references and tend to turn the attention to themselves (I, me, and mine) (0.17). The results support the work of [25], [26], [38]. Their posts contain more negative emotions, sadness and anxiety with a greater emphasis on the present and future. Based on our findings, LIWC can play an effective role for data detection models if added into designing tools. Table 4 summarizes the high-

TABLE 4. Highest correlation results achieved with LIWC features.

LIWC category	Example word	P value
<i>Linguistic Dimensions</i>		
Personal pronoun	I, them, her	0.15*
1st	I, me, mine	0.17***
Person singular		
Negations	no, not, never	0.16**
<i>Psychological processes</i>		
Social Processes	buddy, mate, talk	0.17**
Affective Processes	happy, cry, hate	0.19***
Cognitive Processes	think, know, always	0.14*
<i>Personal Concerns</i>		
Work	job, majors, xerox	0.16***
Money	audit, cash, owe	0.14**
Death	win, success, better	0.11*

Note: **P<0.01, ***P<0.001. All the correlation coefficients meet the P<0.05.

est correlation results achieved with the LIWC dictionary features.

3) PREDICTIVE POWER OF LDA

To calculate the hidden topics extracted from the posts, we built a topic model which functions as a triggering point for depression. LDA requires specifying the number of generated topics. Any change in the parameter may cause a change in the classification accuracy. For this reason, it is necessary to find an appropriate value. As mentioned in section IV-B,

TABLE 5. Example of generated topics with LDA features in depression-related texts.

Topics	Most representative words
Job	boss, boring, broke, company, fired, handjob, jobless, pay, money, quit, stress, left, time, unemployed, unhappy, want, work, year
Depression	always, better, die, depressed, feeling, hate, isolate, long, mind, meaning, guilt, myself, negative, over month, pain, suffer, something, thoughts
Tired	abusive, bullied, burden, can't, hurt, ill, live, loneliness, mentally, myself, neg thought, sleep, started, time, think wrong one, world, wanted
Friends	best, date, dude, chill, care, encourage, happy, help, insecure, relationships, spend, support, no friends, roommate, request, unfollow, worthless
Broke	find, help, heartbroken, emotional, inside lost, love, marriage, often, relationship, rejected, spouse, together, sex, problems, ugly
F***	break, death, drunk, damn, die, done, exercise, long, lonely, pain, kill me, removed, phone, stupid, shit, school, sleep, suck, text, treatment

LDA works well on the validation set limited to 70 topics with the accuracy reaching 70%.

In our study, we select 20 topics with the largest topic proportion from the posts as illustrated in Fig.4. To find the number of topics that allow an optimal classification, we perform several types of testing. We use t-SNE to perform dimensionality reduction and get a two-dimensional coordinates for our data [61]. The t-SNE algorithm has two main hyper-parameters that have an impact on final visualization that is perplexity and number of iterations. We tune them and produce the data representation using the perplexity 50 and 500 iterations. The Fig. 4 shows the topic arrangement by putting related topic clusters close to each other and nicely arranging the documents with different topics.

Table 5 illustrates the words correlated with the particular topics generated from the posts. The topics contain a lexicon of words that is common among the depressed accounts. The topics such as Depression, Broke or Tired reflect the pain, suffering or depressive symptoms of the users with the words related to disclosure. The users feel heartbroken, rejected or lonely with a low self-esteem and self pre-occupation. Additionally, there are other topics which reflect hostility, aggression or relationship with their friends.

B. CLASSIFICATION RESULTS

Our task is to detect depression of each of the users in the chosen data. We start to conduct the execution of the text classifying techniques by using the entire dimension feature

space extracted from the dataset. For baseline features, we use LIWC categories, N-gram probabilities, LDA model and their multiple feature combinations built on the Reddit training data. The aim of combining the distinct NLP techniques is to find out what combination of the features best favors the performance accuracy for depression detection. In this section, we investigate and discuss the degree of accuracy achieved by our experiment. For estimating the presence of depression within the posts, we apply four major classifiers and one artificial neural network classifier defined in IV-C. For implementation of each classifier, we utilize Scikit-learn library for the Python language [42], and use a 10-fold cross-validation to verify the results.

To evaluate the above-mentioned classification techniques, we apply the evaluation metrics, such as accuracy of estimations (Acc.) Eq. (1) and F-score (F1) Eq. (4) consisting of precision (P) and recall (R). It relies on a confusion matrix incorporating the information about each test sample prediction outcome. Accuracy is the rate of correct classification; F1 Eq.(4) score is a harmonic average of the precision and recall; precision estimates how many positively identified samples are correct, recall estimates what proportion of positive samples was correctly identified. The closer both values are, the higher the F1 score is. In the evaluation metrics, we find that there is a number of true positive predictions (TP), true negative predictions (TN), false positive predictions (FP) and false negative predictions (FN) [62]. The most straightforward classifier evaluation score is an accuracy defined as follows:

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \tag{1}$$

$$Precision = \frac{TP}{TP + FP} \tag{2}$$

$$Recall = \frac{TP}{TP + FN} \tag{3}$$

$$F_1 = 2 \cdot \frac{precision \cdot recall}{precision + recall} \tag{4}$$

Table 6 shows the accuracy results of five constructed classification models with the NLP features. Each classified corpus contains an accuracy, F-score, precision and recall result value. While accuracy is a measure used in many studies on depression detection, we also opt for showing precision, recall and F-score which allows us a deeper analysis of the outputs.

Evaluating the performance of single features (LIWC, LDA, unigram/bigram) with the classifiers, we can observe a performance improvement. The best accuracy is achieved with bigrams, especially with SVM learning algorithm resulting in 80% accuracy and 0.79 F1 score; followed by the LIWC feature and RF model (78%, 0.84) and LDA with LR text classifier (77%, 0.83). Although bigram performs well among the single features, considering combined features, it has its significant limitations.

LIWC as a single feature outperforms LDA on predictive power. With RF model it results in the highest accuracy

predictors of depression contained the words related to preoccupation with themselves, feelings of sadness, anxiety, anger, hostility or suicidal thoughts, with a greater emphasis on the present and future.

To measure the signs of depression, we examined the performance of both single feature and combined feature sets using various text classifying methods. Our results show that a higher predictive performance is hidden in proper features selection and their multiple feature combinations. The strength and effectiveness of combined features are demonstrated with the MLP classifier reaching 91% accuracy and 0.93 F1 score achieving the highest performance degree for detecting the presence of depression in Reddit social media conducted in our study.

Additionally, the best feature among the single feature sets is bigram; with SVM classifier it can detect depression with 80% accuracy and 0.79 F1 score. Considering LIWC and LDA features, LIWC outperformed topic models generated by LDA.

Although our experiment shows that the performances of applied methodologies are reasonably good, the absolute values of the metrics indicate that this is a challenging task and worthy of further exploration. We believe this experiment could further underline the infrastructure for new mechanisms applied in different areas of healthcare to estimate depression and related variables. It can be beneficial for the individuals suffering from mental health disorders to be more proactive towards their fast recovery. In our future work, we will try to examine the relationship between the users' personality [64] and their depression-related behaviour reflected in social media.

ACKNOWLEDGMENT

The authors would like to thank Ms. Janka Koperdanova for her full support and editing of the paper.

REFERENCES

- [1] W. H. Organization. (2017). *Depression and Other Common Mental Disorders: Global Health Estimates*. Geneva: World Health Organization. [Online]. Available: <http://www.who.int/en/news-room/fact-sheets/detail/depression>
- [2] M. J. Friedrich, "Depression is the leading cause of disability around the world," *JAMA*, vol. 317, no. 15, p. 1517, Apr. 2017.
- [3] M. Nadeem. (2016). "Identifying depression on twitter." [Online]. Available: <https://arxiv.org/abs/1607.07384>
- [4] S. Paul, S. K. Jandhyala, and T. Basu, "Early detection of signs of anorexia and depression over social media using effective machine learning frameworks," in *Proc. CLEF*, Aug. 2018, pp. 1–9.
- [5] A. Benton, M. Mitchell, and D. Hovy. (2017). "Multi-task learning for mental health using social media text." [Online]. Available: <https://arxiv.org/abs/1712.03538>
- [6] G. Coppersmith, M. Dredze, C. Harman, and K. Hollingshead, "From ADHD to SAD: Analyzing the language of mental health on twitter through self-reported diagnoses," in *Proc. 2nd Workshop Comput. Linguistics Clin. Psychol. Linguistic Signal Clin. Reality*, 2015, pp. 1–10.
- [7] D. Maupomé and M. Meurs, "Using topic extraction on social media content for the early detection of depression," in *Proc. CLEF (Working Notes)*, vol. 2125, Sep. 2018. [Online]. Available: <https://CEUR-WS.org>
- [8] 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.
- [9] D. Preotiu-Pietro et al., "The role of personality, age, and gender in tweeting about mental illness," in *Proc. 2nd Workshop Comput. Linguistics Clin. Psychol. Linguistic Signal Clin. Reality*, 2015, pp. 21–30.
- [10] T. Nguyen, D. Phung, B. Dao, S. Venkatesh, and M. Berk, "Affective and content analysis of online depression communities," *IEEE Trans. Affect. Comput.*, vol. 5, no. 3, pp. 217–226, Jul. 2014.
- [11] H. A. Schwartz et al., "Towards assessing changes in degree of depression through facebook," in *Proc. Workshop Comput. Linguistics Clin. Psychol. Linguistic Signal Clin. Reality*, 2014, pp. 118–125.
- [12] 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.
- [13] J. Wolohan, M. Hiraga, A. Mukherjee, Z. A. Sayyed, and M. Millard, "Detecting linguistic traces of depression in topic-restricted text: Attending to self-stigmatized depression with nlp," in *Proc. 1st Int. Workshop Lang. Cognition Comput. Models*, 2018, pp. 11–21.
- [14] Y. Tyshchenko, "Depression and anxiety detection from blog posts data," *Nature Precis. Sci., Inst. Comput. Sci., Univ. Tartu, Tartu, Estonia*, 2018.
- [15] S. C. Guntuku, D. B. Yaden, M. L. Kern, L. H. Ungar, and J. C. Eichstaedt, "Detecting depression and mental illness on social media: An integrative review," *Current Opinion Behav. Sci.*, vol. 18, pp. 43–49, Dec. 2017.
- [16] R. A. Calvo, D. N. Milne, M. S. Hussain, and H. Christensen, "Natural language processing in mental health applications using non-clinical texts," *Natural Language Eng.*, vol. 23, no. 5, pp. 649–685, 2017.
- [17] Á. Hernández-Castañeda and H. Calvo, "Deceptive text detection using continuous semantic space models," *Intell. Data Anal.*, vol. 21, no. 3, pp. 679–695, Jan. 2017.
- [18] S. Freud, *The Psychopathology of Everyday Life*. London, U.K.: Hogarth, 1901.
- [19] A. T. Beck, *Depression: Clinical, Experim., Theoretical Aspects*. Philadelphia, PA, USA: Univ. Pennsylvania Press, 1967.
- [20] T. Pyszczynski and J. Greenberg, "Self-regulatory perseveration and the depressive self-focusing style: A self-awareness theory of reactive depression," *Psychol. Bull.*, vol. 102, no. 1, p. 122, Jul. 1987.
- [21] E. Durkheim and A. Suicide, *A Study in Sociology*. Abingdon, U.K.: Routledge, 1952.
- [22] A. G. Reece, A. J. Reagan, K. L. Lix, P. S. Dodds, C. M. Danforth, and E. J. Langer, "Forecasting the onset and course of mental illness with twitter data," *Sci. Rep.*, vol. 7, no. 1, p. 13006, Oct. 2017.
- [23] S. Tsugawa, Y. Kikuchi, F. Kishino, K. Nakajima, Y. Itoh, and H. Ohsaki, "Recognizing depression from twitter activity," in *Proc. CHI*, Apr. 2015, pp. 3187–3196.
- [24] J. C. Eichstaedt et al., "Facebook language predicts depression in medical records," *Proc. Nat. Acad. Sci.*, vol. c, no. 44, pp. 11203–11208, Oct. 2018.
- [25] S. W. Stirman and J. W. Pennebaker, "Word use in the poetry of suicidal and nonsuicidal poets," *Psychosomatic Med.*, vol. 63, no. 4, pp. 517–522, Jul. 2001.
- [26] S. Rude, E.-M. Gortner, and J. Pennebaker, "Language use of depressed and depression-vulnerable college students," *Cognition Emotion*, vol. 18, no. 8, pp. 1121–1133, Dec. 2004.
- [27] J. Zinken, K. Zinken, J. C. Wilson, L. Butler, and T. Skinner, "Analysis of syntax and word use to predict successful participation in guided self-help for anxiety and depression," *Psychiatry Res.*, vol. 179, no. 2, pp. 181–186, Sep. 2010.
- [28] T. S. PortalStatistics and Studies. (2019). *Social Media Usage Worldwide*. [Online]. Available: <https://www.statista.com/statistics/272014/global-social-networks-ranked-by-number-of-users/>
- [29] M. De Choudhury, S. Counts, and E. Horvitz, "Predicting postpartum changes in emotion and behavior via social media," in *Proc. SIGCHI Conf. Hum. Factors Comput. Syst.*, Apr. 2013, pp. 3267–3276.
- [30] M. A. Moreno et al., "Feeling bad on Facebook: Depression disclosures by college students on a social networking site," *Depression Anxiety*, vol. 28, no. 6, pp. 447–455, Jun. 2011.
- [31] M. D. Choudhury, E. Kiciman, M. Dredze, G. Coppersmith, and M. Kumar, "Discovering shifts to suicidal ideation from mental health content in social media," in *Proc. CHI Conf. Hum. Factors Comput. Syst.*, May 2016, pp. 2098–2110.
- [32] S. Bagroy, P. Kumaraguru, and M. De Choudhury, "A social media based index of mental well-being in college campuses," in *Proc. CHI Conf. Hum. Factors Comput. Syst.*, May 2017, pp. 1634–1646.
- [33] D. E. Losada, F. Crestani, and J. Parapar, "eRISK 2017: CLEF lab on early risk prediction on the internet: Experimental foundations," in *Proc. Int. Conf. Cross-Lang. Eval. Eur. Lang.*, 2017, pp. 346–360.

- [34] H. Almeida, A. Briand, and M.-J. Meurs, "Detecting early risk of depression from social media user-generated content," in *Proc. CLEF*, 2017, pp. 1–10.
- [35] M. Trotzek, 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.*, to be published.
- [36] N. Masuda, I. Kurahashi, and H. Onari, "Suicide ideation of individuals in online social networks," *PLoS One*, vol. 8, no. 4, Apr. 2013, Art. no. e2262.
- [37] A. Li, X. Huang, B. Hao, and B. O'Dea, H. Christensen, T. Zhu, and A. F. Jorm, "Attitudes towards suicide attempts broadcast on social media: An exploratory study of chinese microblogs," *PeerJ*, vol. 8, no. 3, p. e1209, Sep. 2015.
- [38] N. Ramírez-Esparza, C. K. Chung, E. Kacewicz, and J. W. Pennebaker, "The psychology of word use in depression forums in english and in spanish: Texting two text analytic approaches," in *Proc. ICWSM*, Mar. 2008, pp. 1–10.
- [39] I. Pirina and C. C. Çoltékin, "Identifying depression on reddit: The effect of training data," in *Proc. EMNLP Workshop*, 2018, pp. 9–12.
- [40] G. Coppersmith, M. Dredze, C. Harman, K. Hollingshead, and M. Mitchell, "Clpsych 2015 shared task: Depression and ptsd on twitter," in *Proc. 2nd Workshop Comput. Linguistics Clin. Psychol. Linguistic Signal Clin. Reality*, 2015, pp. 31–39.
- [41] G. Salton and C. Buckley, "Term-weighting approaches in automatic text retrieval," *Inf. Process. Manage.*, vol. 24, no. 5, pp. 513–523, 1988.
- [42] F. Pedregosa *et al.*, "Scikit-learn: Machine learning in Python," *J. Mach. Learn. Res.*, vol. 12, pp. 2825–2830, Oct. 2011.
- [43] C. D. Manning and H. Schütze, *Foundations of Statistical Natural Language Processing*, vol. 999. Cambridge, MA, USA: MIT Press, 1999.
- [44] J. W. Pennebaker, R. J. Booth, R. L. Boyd, and M. E. Francis, "Linguistic Inquiry and Word Count: LIWC2015," Pennebaker Conglomerates, Austin, TX, USA, 2015. [Online]. Available: www.LIWC.net
- [45] C. Chung and J. W. Pennebaker, "The psychological functions of function words," *Social Commun.*, vol. 1, pp. 343–359, Sep. 2007.
- [46] M. Park, D. W. McDonald, and M. Cha, "Perception differences between the depressed and non-depressed users in twitter," in *Proc. 7th Int. AAAI Conf. Weblogs Social Media*, Jun. 2013, pp. 45–65.
- [47] C. M. Homan, N. Lu, X. Tu, M. C. Lytle, and V. Silenzio, "Social structure and depression in trevorspace," in *Proc. 17th ACM Conf. Comput. Supported Cooperat. Work Social Comput.*, Feb. 2014, pp. 615–625.
- [48] G. Coppersmith, M. Dredze, and C. Harman, "Quantifying mental health signals in twitter," in *Proc. Workshop Comput. Linguistics Clin. Psychol. Linguistic Signal Clin. Reality*, 2014, pp. 51–60.
- [49] P. Resnik, A. Garron, and R. Resnik, "Using topic modeling to improve prediction of neuroticism and depression in college students," in *Proc. Conf. Empirical Methods Natural Lang. Process.*, Sep. 2013, pp. 1348–1353.
- [50] D. M. Blei, A. Y. Ng, and M. I. Jordan, "Latent Dirichlet allocation," *J. Mach. Learn. Res.*, vol. 3, pp. 993–1022, Mar. 2003.
- [51] A. K. McCallum. (2002). *MALLET: A Machine Learning for Language Toolkit*. [Online]. Available: <http://mallet.cs.umass.edu>
- [52] S. L. Gortmaker, "Theory and methods—applied logistic regression by david w. hosmer jr and stanley lemeshow," *Contemp. Sociol.*, vol. 23, no. 1, p. 159, Jun. 1994.
- [53] A. Agresti, *An Introduction to Categorical Data Analysis*. Hoboken, NJ, USA: Wiley, 2018.
- [54] W. S. Noble, "What is a support vector machine?," *Nature Biotechnol.*, vol. 24, no. 12, p. 1565, May 2006.
- [55] B. Xu, Y. Ye, and L. Nie, "An improved random forest classifier for image classification," in *Proc. IEEE Int. Conf. Inf. Autom.*, Jun. 2012, pp. 795–800.
- [56] Y. Freund, R. Schapire, and N. Abe, "A short introduction to boosting," *J.-Jpn. Soc. Artif. Intell.*, vol. 14, nos. 771–780, p. 1612, 1999.
- [57] R. E. Schapire, Y. Singer, and A. Singhal, "Boosting and rocchio applied to text filtering," *SIGIR*, vol. 98, pp. 215–223, Aug. 1998.
- [58] H. Ramchoun, M. A. J. Idrissi, Y. Ghanou, and M. Ettaouil, "Multilayer perceptron: Architecture optimization and training," *IJIMAI*, vol. 4, no. 1, pp. 26–30, 2016.
- [59] G. Cybenko, "Approximation by superpositions of a sigmoidal function," *Math. Control, Signals Syst.*, vol. 2, no. 4, pp. 303–314, 1989.
- [60] J. S. Buyukdura, S. M. McClintock, and P. E. Croarkin, "Psychomotor retardation in depression: Biological underpinnings, measurement, and treatment," *Progr. Neuro-Psychopharmacol. Biol. Psychiatry*, vol. 35, no. 2, pp. 395–409, Mar. 2011.
- [61] L. van der Maaten and G. E. Hinton, "Visualizing data using T-SNE," *J. Mach. Learn. Res.*, vol. 9, pp. 2579–2605, Nov. 2008.
- [62] T. Basu and C. Murthy, "A feature selection method for improved document classification," in *Advanced Data Mining and Applications*. New York, NY, USA: Springer, 2012, pp. 296–305.
- [63] X. Huang, T. Liu, A. Li, Z. Chen, and T. Zhu, "Using linguistic features to estimate suicide probability of chinese microblog users," in *Proc. HCC*, Sep. 2014, pp. 145–156.
- [64] M. M. Tadesse, H. Lin, B. Xu, and L. Yang, "Personality predictions based on user behavior on the facebook social media platform," *IEEE Access*, vol. 6, pp. 61959–61969, 2018.



MICHAEL M. TADESSE received the B.S. degree in management information systems (MIS) from Unity University College, Addis Ababa, Ethiopia, in 2005, and the M.S. degree in software engineering from Chongqing University, in 2011. He is currently pursuing the Ph.D. degree with the School of Computer Science and Technology, Dalian University of Technology. His current research interests include personality prediction, recommendation system information retrieval, and natural language processing.



HONGFEI LIN received the B.Sc. degree from Northeastern Normal University, in 1983, the M.Sc. degree from the Dalian University of Technology, in 1992, and the Ph.D. degree from Northeastern University, in 2000. He is currently a Professor with the School of Computer Science and Technology, Dalian University of Technology. He is also the Director of the Information Retrieval Laboratory, Dalian University of Technology. He has published over 100 research papers

in various journals, conferences, and books. His research interests include information retrieval, text mining, natural language processing, and effective computing. In recent years, he has focused on text mining for biomedical literatures, biomedical hypothesis generation, information extraction from a huge biomedical resources, learning to rank, sentimental analysis, and opinion mining. His research projects are funded by the National Natural Science Foundation of China and the National High-Tech Development Plan.



BO XU received the B.Sc. degrees from the Dalian University of Technology, China, in 2011, where he is currently pursuing the Ph.D. degree with the School of Computer Science and Technology. His current research interests include information retrieval, learning to rank, natural language processing, and biomedical text mining.



LIANG YANG received the B.Sc. degree from the Dalian University of Technology, China, in 2009, and the Ph.D. degree, in 2017, where he is currently a Lecturer with the School of Computer Science and Technology, Dalian University of Technology. His current research interests include sentiment analysis, natural language processing, and text mining.