## CHAPTER 1                                           INTRODUCTION

The proliferations of internet and communication technologies, especially the online social networks have rejuvenated how people interact and communicate with each other electronically. The applications such as Instagram, Twitter, Instagram and alike not only host the written and multimedia contents but also offer their users to express their feelings, emotions, and sentiments about a topic, subject or an issue online. On one hand, this is great for users of social networking site to contribute and respond to any topic online; on the other hand, it creates opportunities for people working in the health sector to get insight of what might be happening at mental state of someone who reacted to a topic in a specific manner openly and freely. To provide such insight, machine learning techniques could potentially offer some unique features that can assist in examining the unique patterns hidden in online communication and process them to reveal the mental state (such as 'happiness', 'sadness', 'anger', 'anxiety', depression) among social networks' users. Moreover, there is growing body of literature addressing the role of social networks on the structure of social relationships such as breakup relationship, mental illness ('depression', 'anxiety', 'bipolar' etc.), smoking and drinking relapse, sexual harassment and for suicide ideation . In this study, we aim to analyze Instagram data to detect any factors that may reflect the depression of relevant Instagram's users. Various machine learning techniques are employed for such purpose. Considering the key objective of this study, the following are subsequent research challenges addressed in paper. Define what depression is and what are the common factors contributing toward depression. What are the factors to look for depression detection in Instagram comments? How to extract these factors from Instagram comments? What is the relationship between these factors and attitudes toward depression? When is the most influential time to communicate within depressive Indicative Instagram user? What are the most influential machine learning techniques for detection of depression in Instagram comments? In the context of above-mentioned challenges, we analyze depression from Instagram users' data. As users express their feeling as a post or comments in the Instagram platform, sometimes their posts and comments refer to as emotional state such as 'joy', 'sadness', 'fear', 'anger', or 'surprise'.

We analyze various features of Instagram comments by collecting data through an effective method of machine learning classification techniques and to make overall judgements regarding their various parts. In this study, we used publicly available Instagram data (from bipolar, depression and anxiety Instagram page) containing users' comments. Once we access the data, it was cleaned from any inconsistency and then analyzed by a software

application called LIWC. In this study, we examine various linguistic cues which help to detect emotion cause events: the position of cause event and experiencer relative to the emotion keyword: emotional process like positive emotion (e.g. 'happy', 'love', 'nice'), negative emotion (e.g. 'worthless', 'loser', 'hurt', 'ugly', 'nasty'), sadness (e.g. 'worry', 'crying', 'grief', 'sad'), anger (e.g. 'stop', 'shit', 'hate', 'kill', 'annoyed') and anxiety (e.g. 'worried', 'fearful'). A temporal process like present focus (e.g., 'today', 'is', 'now'), past focus (e.g., 'ago', 'did', 'talked') and future focus (e.g., 'shall', 'may', 'will', 'soon'). Linguistic words like articles (e.g. 'a', 'an', 'the'), prepositions (e.g. 'for', 'in', 'of', 'to', 'with', 'above'), auxiliary verbs (e.g. 'do', 'have', 'am', 'will'), conjunctions (e.g. 'and', 'but', 'whereas'), personal pronoun (e.g. 'I', 'them', 'her', 'him'), impersonal pronouns (e.g. 'it', 'it's', 'those'), verbs (e.g. 'go', 'good') and negation (e.g. 'deny', 'dishonest', 'no', 'not', 'never'). the main contributions of this paper are listed as follows: First, we synthesized the literature on various emotion detection techniques to detect depression. Second, we designated four features for our specific research problem and elaborate on the lesson learned from using each type. Third, our experiments are carried out on datasets of Instagram user comments. Fourth, we suggest machine learning techniques to utilize all factors and maintain robustness. Finally, our work also shows the importance of depression detection for mental disorder detection. the remainder of the paper is organized as follows: "Related work" presents the related work of detecting depression analysis of social network data. Methodology is explained in the third section. The experimental analysis is presented in the fourth section, and its discussion in the fifth section. Finally, the conclusion and future work are provided in the last section

## 1.1 OVERVIEW

Purpose: Social networks have been developed as a great point for its users to communicate with their interested friends and share their opinions, photos, and videos reflecting their moods, feelings, and sentiments. This creates an opportunity to analyze social network data for user's feelings and sentiments to investigate their moods and attitudes when they are communicating via these online tools. Methods: Although diagnosis of depression using social networks data has picked an established position globally, there are several dimensions that are yet to be detected. In this study, we aim to perform depression analysis on Instagram data collected from an online public source. To investigate the effect of depression detection, we propose machine learning technique as an efficient and scalable method. Results: We report an implementation of the proposed method. We have evaluated the efficiency of our proposed method using a set of various psycholinguistic features. We show that our

proposed method can significantly improve the accuracy and classification error rate. In addition, the result shows that in different experiments Decision Tree (DT) gives the highest accuracy than other ML approaches to find the depression. Conclusions: Machine learning techniques identify high quality solutions of mental health problems among Instagram users.

## 1.2  MOTIVATION

Depression is one of the most common mental disorders which not only hampers an individual's daily function, but also causes a variety of serious social problems such as suicide. The depressed individuals may not be aware of their depressed symptoms at all so that they sometimes miss the appropriate time for taking care of depression. To prevent this problem, many researchers used social media to identify depressed users by analyzing the differences in language use. We investigate different aspects of posts by depressed users through four feature networks, which will help researchers to effectively screen social media posts to find useful evidence for depressive symptoms with respect to psychological theories. most previous studies could not explain the detection results adequately with respect to relevant theories in the field, making it difficult to conduct a more detailed analysis for further processes such as diagnosis and prevention.

## 1.3  SCOPE OF THE WORK

This project applies data mining techniques to psychology, specifically the field of depression, to detect depressed users in social network services (SNS). The expansion of data mining to psychology is of great technical and social significance. It is proved that the proposed model in this paper could effectively help for detecting depressed ones and preventing suicide in online social networks. The system promises the following advantages: From the aspect of methodology, data mining techniques is expanded to depression area. Sentiment analysis algorithm, specifically for Chinese micro-blog is proposed for calculating the depression inclination. An association model is established between features abstracted from Micro-blog system and depression inclination. The model also determines the principal features which affecting depression detection significantly. An application in Sina Micro-blog is developed for monitoring the users' mental health in SNS. The basic idea of this paper could be explicitly extended to other language scenarios.

## 1.4  PROBLEM STATEMENT

Depression is the world's fourth largest disease and will be in the second place in 2020 according to World Health Organization statistics [1]. It has long been a hot research subject in psychology. With the rise of social networks, many psychologists in depression studies are now turning their sights to web medias from traditional case studies. They detect depressed users with psychological diagnostic criteria and observe the online features of them [2], [3]. Most of the observations are about the behaviors of each individual, and little attention is paid to the interactions between users. Taking the typical features derived from psychologist's research, Wang proposed a depression detection model to classify the users in social networks to be depressed or not [4]. These features are obtained through user profile and their microblog content, such as user activeness, sentiment analysis of micro-blog content and the percent of original blogs. Apparently, the features in Wang's model are all about the node itself, and no linkage feature is considered, so in this paper we call it the classifier with node features considered only (NFO)..

## 1.5  OBJECTIVES

In this model, we aim to analyze Instagram data to detect any factors that may reflect the depression of relevant Instagram's users. Various machine learning techniques are employed for such purpose. Considering the key objective of this study, the following are subsequent research objectives addressed in project.

- Data set exploration.
- Data set preparation.
- Building ground truth dataset.
- Feature extraction.
- Measuring depressive behavior.
- Experimental analysis.

**CHAPTER 2**                                                     **LITERATURE SURVEY**

In article [1], Datasets originating from social networks are valuable to many fields such as sociology and psychology. But the supports from technical perspective are far from enough, and specific approaches are urgently in need. This paper applies data mining to psychology area for detecting depressed users in social network services. Firstly, a sentiment analysis method is proposed utilizing vocabulary and man-made rules to calculate the depression inclination of each micro-blog. Secondly, a depression detection model is constructed based on the proposed method and 10 features of depressed users derived from psychological research. Then 180 users and 3 kinds of classifiers are used to verify the model, whose precisions are all around 80%. Also, the significance of each feature is analyzed. Lastly, an application is developed within the proposed model for mental health monitoring online. This study is supported by some psychologists and facilitates them in data-centric aspect in turn. A camera is used to capture the images of the faces or to capture the real-time video. Optical devices such as the camera or video recorder are used to accomplish this task. The students face images to be recognized are fed to the image processing block where it performs preprocessing, face detection, and face recognition tasks. Preprocessing includes tasks such as cropping of image and enhancement procedures. These processed images are fed to the face recognition algorithm. These database images are then compared with the Realtime recognized faces to identity the student.

Article [2] is a proposed work on social networks contain a tremendous amount of node and linkage data, providing unprecedented opportunities for a wide variety of fields? As the world's fourth largest disease, depression has become one of the most significant research subjects. Previously, a depression classifier has been proposed to classify the users in online social networks to be depressed or not, however, the classifier takes only node features into account and neglects the influence of linkages. This paper proposes an improved model to calculate the probability of a user being depressed, which is based on both node and linkage features. The linkage features are measured in two aspects: tie strength and interaction content analysis. Moreover, the propagation rule of depression is considered for improving the prediction accuracy. Finally, our experiments on the data derived from Sina Micro-blog shows that the highest accuracy of the improved model is 60%, increasing by algorithm compared to the classifier with node features considered only. In this paper, it is well proved that adding linkage features analysis performs much better than node features analysis only. It also implies that tie strength and interaction content have different effects on depression probability

estimation. Although this model is proposed for depression detection, the basic idea of linkage features analysis could be explicitly used in a wide scenario.

Article [3] is presented with social networks have been developed as a great point for its users to communicate with their interested friends and share their opinions, photos, and videos reflecting their moods, feelings, and sentiments. This creates an opportunity to analyze social network data for user's feelings and sentiments to investigate their moods and attitudes when they are communicating via these online tools. Although diagnosis of depression using social networks data has picked an established position globally, there are several dimensions that are yet to be detected. In this study, we aim to perform depression analysis on Facebook data collected from an online public source. To investigate the effect of depression detection, we propose machine learning technique as an efficient and scalable method. We report an implementation of the proposed method. We have evaluated the efficiency of our proposed method using a set of various psycholinguistic features. We show that our proposed method can significantly improve the accuracy and classification error rate. In addition, the result shows that in different experiments Decision Tree (DT) gives the highest accuracy than other ML approaches to find the depression.

In article[4], the proposed system although depression is one of the most common mental disorders, the depressed individuals may not be aware of their symptoms at all so that they sometimes miss the appropriate time for treatment. To prevent this problem, many researchers investigated social media to figure out depressed individuals by analyzing the differences in language use. While they have recently achieved reasonable performance in detecting depression, especially using deep learning methods, such methods still do not provide a clear way to explain why certain individuals have been detected as depressed. To address this issue, we propose Feature Attention Network (FAN), inspired by the process of diagnosing depression by an expert who has background knowledge about depression. We evaluate the performance of our model on a large-scale general forum (Reddit Self-reported Depression Diagnosis) dataset. Experimental results demonstrate that FAN shows good performance with high interpretability despite a smaller number of posts in training data. We investigate different aspects of posts by depressed users through four feature networks built upon psychological studies, which will help researchers to investigate social media posts to find useful evidence for depressive symptoms.

**CHAPTER 3**                                                    **METHODOLOGY**

In this study, we first focused on four types of factors such as emotional process, temporal process, linguistic style, and all (emotional, temporal, linguistic style) features together for the detection and processing of depressive data received as Instagram posts. We then apply supervised machine learning approaches to study each factor types independently.

**Data set exploration**

We worked on Instagram users' comments for depressive behavioral exploration and detection. We collected data from the social network . Preparing of social network data, in particular Instagram user's comments is one of the primary challenges which bear information on whether they could contain depression bearing content. To tackle this issue, we use Capture for collecting data from Instagram. For qualitative data analysis, Capture is a powerful tool in the world today. It is intended to enable to arrange, break down and discover knowledge in unstructured data like open-ended survey responses, social media, interviews, articles, and web content. Furthermore, it gives a place to arrange and deal with material to discover knowledge in a more proficient way.

**Data set preparation**

After collecting the raw data from Facebook, it was analyzed by using LIWC Software. LIWC is the heart of the text analysis strategy and can process text on a line by line. Our primary dataset contains total 21 columns where 13 columns represent the linguistic style (articles, prepositions, auxiliary verbs, conjunctions, personal pronoun, impersonal pronouns, verbs, negation etc.) information, 5 columns represent the emotional (positive, negative, sad, anger and anxiety) information, 3 columns represent the temporal process (past, present and future) information and each column gives the individual information's about depressive behavior

**Building ground truth dataset**

This section discusses the process employed to construct our dataset with ground truth label information (on whether the comments are depression indicative). The Facebook data containing users' comments were divided into two sets (a) for the positive (YES) class (depression indicative comments) and (b) for the negative (NO) class (non-depression indicative comments.
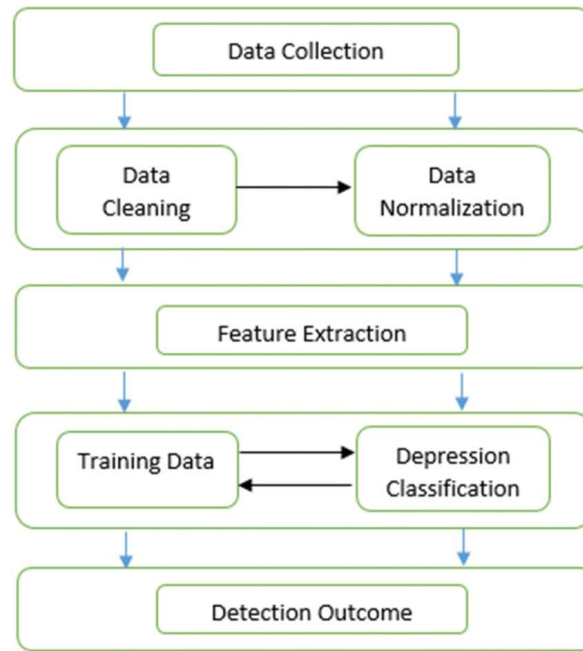
**Fig. 3.1** A methodological overview of Facebook data analysis for depression analysis

**Feature extraction**

To describe and demonstrate amongst depressive and non-depressive posts, we extract the different features in view of psycholinguistic measurements from the user's post. It is clarified briefly as follows: Psycholinguistic features LIWC is a psycholinguistic vocabulary package made by psychological analysts to perceive the different affective, intellectual, and etymological parts lies on user's verbal or written correspondence. It returns more than 70 different factors with higher level of psycholinguistic features, for example

• Psychological process—affective process, social process, cognitive process, perceptual process, biological process, drives, time orientations, relativity, personal concerns •Linguistic process—word count, word/sentence, pronoun, personal pronoun, articles, prepositions, auxiliary verbs, adverbs, conjunctions, Negations

• Other's grammar—verbs, adjectives, comparisons, interrogatives, number, quantifiers.

These higher-level categories are also divided into subcategories such as

• Biological processes—sexual, body, ingestion, and health

• Affective processes—anxiety, anger, sadness, positive emotion, negative emotion

• Time orientations—present, past, future

• Social processes—family, friends, male, female

• Perceptual processes—see, hear, feel

**Measuring depressive behavior**

We presented a set of attributes like emotional process, temporal process, and linguistic style that can be used to characterize the depressive behaviors of users. Our dataset consists of five emotional variables (positive, negative, sad, anger, anxiety), three temporal categories (present focus, past focus, and future focus), and 9 standard linguistic dimensions (e.g., articles, prepositions, auxiliary verb, adverbs, conjunctions, pronoun, verbs, and negations.

**Temporal process**

Generally, temporal process word provides information about past focus category, present focus category and future focus category of how people are referencing each other and their degree of emotionality.

**Linguistic process**

Linguistics process is one of the largest parts of LIWC psycholinguistics vocabulary package. It was intended to quantify word use in mentally significant classifications. Also, it has been effectively used to recognize connections between people in social co-operations, including relative status, trickiness, and the nature of close relationship. So, in our study we use nine specific linguistics features (articles, prepositions, auxiliary verbs, adverbs, conjunctions, personal pronoun, impersonal pronouns, verbs, and negations) to characterize user comments for our experimental analysis.

**Classification model**

This stage constructs prediction model for depression post/comments recognition, by considering the psycholinguistic features as input. the task of a classifier f is to find the corresponding label for each posts/comment and k-Nearest Neighbor (kNN). It is a non-probabilistic linear binary classifier that analyzes data for classification or anomaly detection. It builds a hyperplane into high dimensional feature space and finds a hyperplane that isolates the data into two classes with the biggest separation to the closest training data purpose of any class methods use multiple learning algorithms of decision tree for better predictive

performance. K-Nearest Neighbor (KNN) K-Nearest Neighbor (KNN) is a non-parametric approach use to discover the distances from point of interest to points in training set.

**Experimental analysis**

In this study, we examine the execution of various classifiers for depression detection in a shorter time.

**Data analysis**

We applied two major classifiers:, K-Nearest Neighbors (KNN), and Random Forest Classifier.. To comprehend the significance of different feature types, we applied four classifiers' techniques each utilizing: emotional process, linguistic style, temporal process, and all features. the results of the analysis are reported in Fig 2 and that suggests Random Forest Classifier as best performing model. Although KNN gives the near precision, but Random Forest Classifier gives the highest result for recall and F-measure relating to the class of depression indicative comments of Instagram user. Similarly, for linguistic style Random Forest Classifier gives the highest result for precision, recall and F-measure
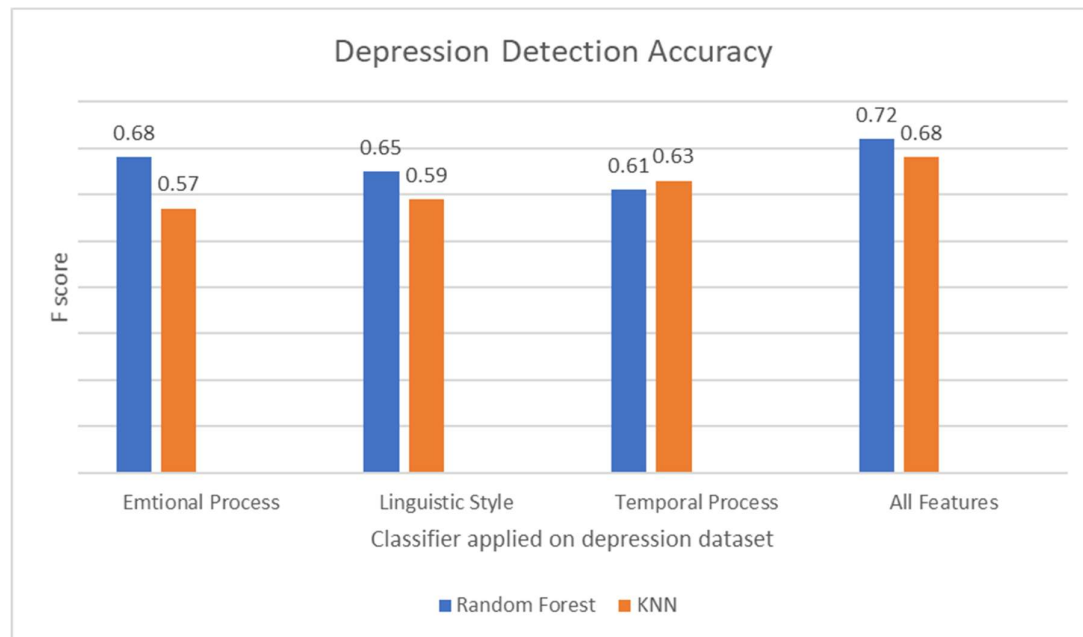


**Fig 3.2** Depression Detection Accuracy

**CHAPTER 4**                                            **SENTIMENT ANALYSIS TECHNIQUES**

As the cardinal symptom of depression is severe negative emotions and lack of positive emotions, sentiment analysis is the most important step in depression detection. Sentiment analysis aims at mining users' opinions and sentiment polarity from the texts they posted [7]. Recently many progresses have been made in sentiment analysis on Twitter data. This research includes two aspects: • Subject-independent analysis, namely judging the polarity of the tweets without considering if it is relevant to a subject [8-10]. The main approaches are based on hashtags, smileys, and some abstract features. • Subject-dependent analysis, namely judging the polarity of the tweets based on the given subject [11-12]. The sentiments of the tweets as positive, negative, or neutral in [11], according to not only the abstract features but also the target-dependent features, which refers to the comments on the target itself and the related things, which are defined as extended targets. Sentiment analysis research on Chinese text is still in its starting stage [13]. The highest accuracy of polarity discrimination on Chinese text is only 59.27% in the latest NTCIR evaluation [14]. Little study has been made for solving problems in a specific field, although analysis strategy differs a lot for different fields. For example, depression sufferers tend to think the topic about "death", so this kind of words should be paid special attention to when constructing the vocabulary. Micro-blogs are often written in a colloquial style, which also bring new challenges when instituting the linguistic rules in the proposed method. The problem addressed in this paper is subject-dependent sentiment analysis of micro-blogs. Inspired by the work in literature [11], abstract features and target-dependent features are considered. This study stresses the particularity of depression and micro-blog content, and the whole model is specifically designed based on them. As shown in Fig.1, a sentiment analysis method is firstly proposed utilizing vocabulary and man-made rules to calculate the depressive inclination of each micro-blog in Fig.1 (A). The vocabulary and man-made rules in sentiment analysis method are constructed based on the Chinese syntax rules, the particularity of depression and micro-blogs (section 3). Then as shown in Fig.1 (B), a depression detection model is constructed based on the proposed method and 10 features of depressed users derived from psychological research (section 4). Lastly, the significance of each feature is analyzed, and a simplified model is proposed for the application in Sina Micro-blog.

## Sentiment Analysis of Micro-blog Content

The most direct expression of depressed mood is the users' micro-blog content, so the sentiment analysis method in this section helps to figure out the polarity of each piece of micro-blog, which emphasizes the depression inclination reflected from the content. A vocabulary is constructed based on HowNet [15], and the sentence structure patterns, and calculation rules are derived according to Chinese syntax rules. As described above, the particularity of depression and micro-blogs are paid special attention to the whole process.

## Vocabulary Construction

The most essential particularity of depression and micro-blogs is the use of words. A vocabulary fitting for depression detection is constructed based on How-Net, a comprehensive vocabulary of Chinese words, as shown in Table.

**Table 4.1** Words in HowNet vocabulary

| 2 Item | | Number | Example |
|---|---|---|---|
| Emotion Words | Positive | 4566 | Pretty, love, like, happy, good |
| | Negative | 4370 | Ugly, sad, depressed, unhappy, bad |
| Degree Modifiers | | 219 | Most(2), over(1.75), very(1.5), more(1), -ish(0.75), insufficient(0.5) |

HowNet contains most of the popular emotion words and degree modifiers. The weights of degree modifiers are quantified into six levels according to their intensities. HowNet is designed for general sentiment analysis. To make it fit for depressed inclination calculation, several adjustments are made as follows:

1. Emotion words, cyberspeaks, modal particles and negative words are added: 1) Depressed users tend to use more emotion words, especially negative emotion words [4], some of which are even only for them. For example, "bye" is a neutral word for normal people, but

it's a typical negative one for depressed users. So, these typical emotion words for depression are added. 2) Considering that cyberspeaks are in prevalent in Internet, they are playing an important role in micro-blogs. Therefore, these words are also added, for example, "smilence", which means "smile silently", into the vocabulary. 3) As micro-blogs are often written in a colloquial style, modal particles often occur in micro-blogs to express feelings directly, such as "ha-ha" and "a-ha", so these modal particles are added into the vocabulary too. 4) Besides degree modifiers, negative words could also modify the expressions, which do not exist in HowNet. Negative words totally reverse the meaning of the expression, such as "not", "never". The negative words selected from a dictionary are imported as a new item into the vocabulary. 2. The part of speech of each word are recognized: The proposed calculation rules is derived from Chinese syntax rules, which are defined by parts of speech. So, the part of speech of each word are recognized and imported as an attribute into the vocabulary. Finally, the vocabulary is constructed with three items as shown in Table 2. 1210 emotion words and 36 negative words are added. Each word holds its own part of speech and weight. Degree modifiers are inherited from HowNet.

## Linguistic Rules Construction

The meaning of a sentence could not be decided only by the words it uses, but also by the order of words, named the structure of the sentence. so, the structure of sentence should be considered in the process of polarity calculation. The structure of sentences could be described as the linguistic rules, which reflects the complexity of Chinese language in one aspect. In this section, linguistic rules based on the proposed vocabulary is constructed by taking the colloquial style of micro-blog into account. According to Chinese syntax rules, the proposed linguistic rules are derived as shown in Table 3. In the rules, Sentence structure patterns are recognized with different items of words in the vocabulary, and each pattern has its own calculation rule based on the weight of each word. If the sentence is recognized as "Partial Negative Structure", a coefficient should be multiplied for precise result, which is set as -0.5 in Table 3. How to calculate the polarity of a given sentence according to the rules will be introduced in the next section.

## Procedure of the Proposed Method

Within the preparation of vocabulary and linguistic rules construction, the proposed method contains 3 main steps

**Sentence Segmentation and Word Segmentation**: A piece of micro-blog allows 140 characters at most, so it may contain several sub-sentences. Punctuations are taken as symbols to segment sentences. ICTCLAS Chinese word segmentation systems, the most popular one throughout the world, is applied for segmenting word and labeling part of speech of each word.

• Each micro-blog is regarded as a sentence S, and each S is a sequence of N sub-sentences denoted by S= {s1,s2,...sn} , where sn is the nth sub-sentence.

• A sub-sentence si is a collection of M words denoted by si={w1(sp1),w2(sp2) ,...,wM (spM ) }, where wm(spm) is the mth word in the collection, and spm refers to its part of speech.

**Polarity Calculation of Each Sub-sentence**. After being segmented, the polarity of each sub-sentence si could be calculated by structure pattern mining and the corresponding calculation rules. The process is implemented as follows:

| Algorithm 1. Polarity calculation algorithm |
|---|
| **1**: **Sub-sentence si**: I am extremely happy and very glad today. |
| **2: Word segmentation**: I(noun), today(noun), extremely(adverb), happy(adjective), and(adverb), very(adverb), glad(adjective). |
| **3: Keyword extraction in vocabulary**: extremely(adverb)#W DM, happy(adjective)#W EW, very(adverb)#W DM, glad(adjective)#W EW |
| **4: Structure pattern mining**: W EW+W EW: happy(adjective)#W EW+ glad(adjective)#W EW; W DM+W EM: extremely(adverb)#W DM+ happy(adjective)#W EW; W DM+W EM: very(adverb)#W DM + glad(adjective)#W EW. |
| **5: Polarity calculation of sub-sentence si**: p(si)= [weight(extremely)×weight(happy)] +[weight(very)×weight(glad)]= [2×(+1)]+[1.5×(+1)]= +3.5. |

## Polarity Calculation of Sentence S

The polarity of sentence S is determined by the polarities and positions of its sub-sentences, as the position of a subsentence si in S can indicate its importance [17]. This is especially noticeable in micro-blog because it enables people to record their immediate feelings at any time. If a micro-blog content is long, the beginning and ending sub-sentences often reflect the writer's feelings more directly, thus more important than those in the middle. Therefore, the higher weights are assigned to sub-sentences at the two ends of the micro-blog as eq(1). With the polarity and position of each sub-sentence, the polarity of S is calculated as eq(2).

$$\lambda(si) = 1/ \min(i, N - i + 1), 1 \leq i \leq N \qquad \text{eq (1)}$$

$$Polarity(S) = N\ i{=}1\ [\lambda(si) \times p(si)] \qquad \text{eq (2)}$$

N is the number of sub-sentences in S, and i is the position of si in S. A positive polarity(S) means the sentence expresses a positive sentiment of the user, and a negative one means opposite. If the polarity equals to zero, then it is objective. The absolute value |polarity(S)| shows the intensity of the sentiment. Research in psychology shows that depressed individuals focus more on negative aspects of their lives. So, the polarity of users' micro-blog contents is an important feature in depression detection in section 4, which is normalized as eq(3), where |S| is the total number of micro-blogs during a given period.

$$NormalizedSentencePolarity = \sum |S|\ Polarity(S)\ /|S| \qquad \text{eq (3)}$$

CHAPTER 5                                      DEPRESSION DETECTION MODEL

When it comes to depression detection, many other features need to be considered. In this section, the work of psychologists is firstly introduced, and then the classifier based on their work is designed for depression detection.

# Psychologists' Work

Psychologists observe the online behaviors of depressed users and discover potential features that could be used to distinguish depressed and non-depressed individuals. All these features mainly come from three dimensions: micro-blog content, interactions, and behaviors [4-6]. Table 4 lists the statistical data of ten features of two depressed and two normal samples in two weeks, in which A and B are the anonymized users. It reveals that most features show significant differences. For example, depressed users tend to use more first-person singular pronouns but less emoticons. However, some features show little influence on these four users, such as times of being forwarded and commented. The proposed model obtained by training data will further illustrate the significance of each feature for depression detection in section 5.

# Model Construction

Taking the achievement of psychologists as background knowledge, their observations need to be converted to parameters that are easily imported into the model as shown in Fig.1 (B). As the calculation of sentence polarity has been discussed in section 3, how to obtain, process and normalize other features will be introduced in this part..

**The Use of First Person Singular and Plural Pronouns**. As the result in [4] shows, depressed users tend to focus on themselves and detach from others. They use more first-person singular pronouns ("I") but less first-person plural pronouns ("We"). So, the use of first-person pronouns is considered in two aspects: • The quantity of first-person pronouns, reflecting their focuses on themselves. • The ratio of first-person singular pronouns to first person plural pronouns as (4), where Qf s and Qfp represent the quantities of first person singular and plural pronouns respectively

For this purpose, all the first person singular and plural pronouns in users' microblogs need to be detected. It requires that the first-person pronouns must be the subject of the sentence. So besides detecting all the first-person words, whether they are the subject of the sentence or not should also be checked. To solve this problem, we choose to check if the word

following them could be used as predicate. If they could, this sentence is considered as the first-person pronoun. Obviously, this method would meet problems when the sentence structure is too complex, for example, some adverbials are following the subject. But as discussed above, micro-blogs are often written in a colloquial style, and complex structures are not frequent, so the proposed method is effective after being tested.

**User Behaviors in Micro-blog**. As the discovery in , frequencies of the posting original blogs could indicate depression levels of the user. So, the percentage of original micro-blogs is taken as one of the features in user behaviors, calculated as (number of original posts)/|S|. It is also found in that the period users post micro-blogs is another indicator of depression level. Table 4 reveals depressive ones tending more active between 0:00-6:00a.m. Therefore, the percentage of micro-blogs posted in this period is calculated as (number of micro-blogs posted during 0:00-6:00a.m.)/|S|, which is used as another feature in user behaviors

**CHAPTER 6**                                                      **EXPERIMENT**

# Data Acquisition and Experiment Result

The proposed model is applied to detect depressed users in Sina Micro-blog, a social network service like Twitter. It is one of the most influential SNS in China. During August 1st -15th, 2012, a group of psychologists made diagnosis on hundreds of volunteers with the means of questionnaire and interviews. They identified 122 depressed sufferers and 346 normal ones. Among them 90 depressed and 90 non-depressed users who use Sina Micro-blog are picked as training dataset. Their information during August 1st -15th are collected through Sina Micro-blog Open Platform API [18]. A total of 6,013 micro-blogs are collected, of which the user who owns the most micro-blogs owns 173, and the least one owns 3 pieces. Over 50,000 sub-sentences are obtained after sentence segmentation. Our experiment is based on these data. After data processing with methods in section 3&4, ten features are obtained for depression detection. Waikato Environment for Knowledge Analysis (Weka), one of the most useful tools for classification [19], is applied to help classify the users into normal or depressed category. To ensure the result being more reliable, three different kinds of classification approaches are employed: Bayes, Trees and Rules [19-20]. The result is obtained with 10-fold cross validation in Fig.2 (A). ROC Area refers to the area under ROC curve, measuring the quality of a classifier. F-measure is the harmonic mean of precision and recall, denoting the accuracy of a classifier comprehensively. The result in Fig.2 (A) reveals the precisions of the proposed model with different classifiers are all around 70%, which is considered acceptable by psychologists to detect depressed users in SNS. Among the incorrect cases, the number of individuals incorrectly classified into normal category and incorrectly into depressed category are approximately equal. Most of these individuals own less than 10 pieces of micro-blogs, indicating that lack of data information would bring error to the proposed model. According to the psychological research, only 3 out of 1000 users online are depressed, so it's very difficult to conduct experiments on large data. However, more than 100 non-depressed familiar friends in Micro-blog are tested with the application in section 6, and more than 85% of them are correctly classified.

## Model Simplification

Besides verifying the effectiveness of the 10 features in the proposed model, the significance of each feature is also studied. Binary logistic regression analysis with Statistical Product and Service Solution (SPSS) is applied to evaluate the significance of each feature in the model [19]. The result is shown in Fig.2 (B). A lower Sig. represents that it is more important. The threshold for Sig. is set as 0.1 and five features are selected for simplifying the model. Experiments show that the precision of the simplified model is declined by less than 5%, however, the bytes of data needed to collect and computing time are significantly reduced. So, the application in section 6 is developed with the simplified model. Furthermore, as Fig.2 (B) shows, the total number of emoticons and original micro-blogs are the most important features, and times of mentioning others and being forwarded are the least important ones. This is a little different from the observation of psychologists in Table 4, which shows times of mentioning others could easily distinguish depressed and non-depressed individuals. It may enlighten psychologists about some further research.

# CHAPTER 7 APPLICATION

Since mental health problem has become the most serious one for modern urban people, the proposed method of sentiment analysis and depression detection model in this paper can be applied in social network services for user mental health state assessment and monitoring. For this purpose, an application in Sina Micro-blog is developed named "Mental Health Testing". The application provides two functions. One is to calculate the polarity of each piece of micro-blog with the sentiment analysis method, which reflects whether the user is optimistic or not. A user is "very optimistic", "a little optimistic" or "pessimistic" according to the total popularity of his latest micro-blogs in a week. The other function is to analyze whether the given user is inclined to be depressed or not with the simplified depression detection model. If a user is tested to be depressed, the application also provides diagnostic messages including some suggestions from the psychologists on active self-regulating strategies.

# CONCLUSION

A model for detecting depressed users in social network based on sentiment analysis is proposed. The sentiment analysis method pays special attention to the characteristics of depression, and ten features are applied in the depression detection model. The precisions obtained from training dataset are all around 60%, and the significance of each feature in the model is also analyzed for model simplification. An application in Micro-blog is developed to test the polarity of micro-blog and the mental state of users, which has helped psychologists detect several potential depressed users in Micro-blog. Although the depression detection model is proposed, the basic idea of the frame especially the sentence structure pattern mining and principle micro-blog features related to depression could be explicitly extended to other language scenarios. In this work, the training data detected by psychologists are widely scattered in social network. It is hard to analyze the relationship between them, so user interaction is paid little attention and simply three parameters are considered. However, homophily is manifest in the group of depression users, which means, the friends of the depression are more likely to be depressive, and different kinds of interactions indicate different results. Thus, the influence of ties between users is contemplated to be studied in the future and a deeper understanding about depression in SNS will be provided. There are still many works to do for improvement. As only 3 out of 1000 users in online social networks are depressed, it's very difficult to conduct experiments on large data. In the future, we hope to detect more appropriate users to further verify our model. With more experiments, we can have a deeper understanding about the influences of each parameter on the model and provide more specific suggestions on how to determine their values in different situations. To make the model more practical, we should understand the depression more intensively. Depression includes many different types, such as manic type, depressed type, psychotic type, and other types. Patients in different types share different characteristics on their behaviors and interaction activities with others. Therefore, a more calibrate analysis of the users' behaviors will lead the model to be more accurate, which calls for a deeper discipline integration of data mining and psychology research in the future.

# REFERENCES

[1] Scott J. *Social network analysis*. *Thousand Oaks*: Sage; 2017.

[2] Serrat O. *Social network analysis. In: Knowledge solutions*. Singapore: Springer; 2017. p. 39–43.

[3] Mikal J, Hurst S, Conway M. *Investigating patient attitudes towards the use of social media data to augment depression diagnosis and treatment: a qualitative study. In: Proceedings of the fourth workshop on computational linguistics and clinical psychology—from linguistic signal to clinical reality*. 2017.

[4] Conway M, O'Connor D. *social media, big data, and mental health: current advances and ethical implications*. Curr Opin Psychol. 2016;9:77–82.

[5] Ofek N, et al. *Sentiment analysis in transcribed utterances. In: PacifcAsia conference on knowledge discovery and data mining*. 2015. Cham: Springer.

[6] Yang Y, et al. *User interest and social influence-based emotion prediction for individuals. In: Proceedings of the 21st ACM international conference on Multimedia*. 2013. New York: ACM.

[7] Tausczik YR, Pennebaker JW. *The psychological meaning of words: LIWC and computerized text analysis methods.* J Lang Soc Psychol. 2010;29(1):24–54.

[8] Pennebaker JW, Francis ME, Booth RJ. *Linguistic inquiry and word count*: LIWC 2001, vol. 71. Mahway: Lawrence Erlbaum Associates; 2001. p. 2001.

[9] Holleran SE. *The early detection of depression from social networking sites*. Tucson: The University of Arizona; 2010.

[10] Greenberg LS. *Emotion-focused therapy of depression. Per Centered Exp Psychother*. 2017;16(1):106–17.

[11] Haberler G. *Prosperity and depression: a theoretical analysis of cyclical movements.* London: Routledge; 2017.

[12] Guntuku SC, et al. *Detecting depression and mental illness on social media: an integrative review*. Curr Opin Behav Sci. 2017; 18:43–9.

[13] De Choudhury M, et al. *Predicting depression via social media*. In: ICWSM, vol. 13. 2013. p. 1–10. 14. De Choudhury M, Counts S, Horvitz E. *Predicting postpartum changes in emotion and behavior via social media. In: Proceedings of the SIGCHI conference on human factors in computing systems*. New York: ACM; 2013. 15. O'Dea B, et al. Detecting suicidality on Twitter. Internet Interv. 2015;2(2):183–8.

[14]     Zhang L, et al. *Using linguistic features to estimate suicide probability of Chinese microblog users. In: International conference on human centered computing*. Berlin: Springer; 2014.

[15]     Aldarwish MM, Ahmad HF. *Predicting depression levels using social media posts. In: 2017 IEEE 13th international Symposium on Autonomous decentralized system (ISADS)*. 2017.

[16]     Zhou J, et al. *Measuring emotion bifurcation points for individuals in social media. In: 2016 49th Hawaii international conference on system sciences (HICSS)*. 2016. Koloa: IEEE.

[17]     Wang X, et al. *A depression detection model based on sentiment analysis in micro-blog social network. In: Trends and applications in knowledge discovery and data mining (PAKDD)*. 2013.

[18]     Nguyen T, et al. *Afective and content analysis of online depression communities. IEEE Trans Afect Comput*. 2014;5(3):217–26.

[19]     Park M, McDonald DW, Cha M. *Perception diferences between the depressed and non-depressed users in Twitter*. In: ICWSM, vol. 9. 2013. p. 217–226.

[20]     Wee J, et al. *The influence of depression and personality on social networking. Comput Hum Behav*. 2017;74:45–52.

[21]     Bachrach Y, et al. *Personality and patterns of Instagram usage. In: Proceedings of the 4th annual ACM web science conference*. 2012. New York: ACM.

[22]     Ortigosa A, Martín JM, Carro RM. *Sentiment analysis in Instagram and its application to e-learning. Comput Hum Behav*. 2014;31:527–41.

[23]     Shen G, et al. *Depression detection via harvesting social media: A multimodal dictionary learning solution. In: Proceeding of the twenty-sixth international joint conference on artificial intelligence (IJCAI-17)*. 2017. p. 3838–3844.

[24]     https://github.com/ranju12345/Depression-Anxiety-Instagram-pageComments-Text.

[25]     Bazeley P, Jackson K. *Qualitative data analysis with NVivo*. London: Sage; 2013.

[26]     AlYahmady HH, Alabri SS. *Using NVivo for data analysis in qualitative research. Int Interdiscip J Educ*. 2013;2(2):181–6.

[27]     Bandara, W. *Using Nvivo as a research management tool: a case narrative. In: Quality and impact of qualitative research: proceedings of the 3rd international conference on qualitative research in IT & IT in Qualitative Research*. 2006.