

Towards employing social media for studying mental health

by

Amir Hossein Yazdavar

B.S., Shiraz University, 2011

M.S., University Technology Malaysia, 2013

AN ABSTRACT OF A DISSERTATION

submitted in partial fulfillment of the
requirements for the degree

DOCTOR OF PHILOSOPHY

Department of Computer Science
Carl R. Ice College of Engineering

KANSAS STATE UNIVERSITY
Manhattan, Kansas

2021

Abstract

With the ubiquity of social media platforms, millions of people are now sharing their online persona by expressing their thoughts, moods, emotions, and even their daily struggles with mental health on social media. Unlike traditional observational cohort studies conducted through questionnaires and self-reported surveys, we explore the reliable detection of clinical depression from social media obtained unobtrusively.

First, we developed a semi-supervised statistical model to evaluate how the duration of these symptoms and their expression on Twitter align with the medical findings reported via the PHQ-9. Based on the analysis of tweets crawled from users, we demonstrate the potential of detecting clinical depression symptoms which emulate the PHQ-9 questionnaire clinicians use today.

Over the course of this dissertation, we examine and exploit multi-modal big (social) data to discern depressive behaviors using a wide variety of features including individual-level demographics. By developing a multi-modal framework and employing statistical techniques to fuse heterogeneous sets of features obtained through the processing of visual, textual, and user interaction data, we significantly enhance the current state-of-the-art approaches for identifying depressed individuals on social media as well as facilitate demographic inferences from social media. Besides providing insights into the relationship between demographics and mental health, our research assists in the design of a new breed of demographic-aware health interventions.

Altogether, these research topics, resulted in a framework, that when executed, will assist in identifying community-level risk and protective factors associated with the diagnosis and treatment of depression that could be an efficient means of studying patterns of access and utilization of mental health services to inform interventions.

Towards employing social media for studying mental health

by

Amir Hossein Yazdavar

B.S., Shiraz University, 2011

M.S., University Technology Malaysia, 2013

A DISSERTATION

submitted in partial fulfillment of the
requirements for the degree

DOCTOR OF PHILOSOPHY

Department of Computer Science
Carl R. Ice College of Engineering

KANSAS STATE UNIVERSITY
Manhattan, Kansas

2021

Approved by:

Major Professor
Pascal Hitzler

Copyright

© Amir Hossein Yazdavar 2021.

Abstract

With the ubiquity of social media platforms, millions of people are now sharing their online persona by expressing their thoughts, moods, emotions, and even their daily struggles with mental health on social media. Unlike traditional observational cohort studies conducted through questionnaires and self-reported surveys, we explore the reliable detection of clinical depression from social media obtained unobtrusively.

First, we developed a semi-supervised statistical model to evaluate how the duration of these symptoms and their expression on Twitter align with the medical findings reported via the PHQ-9. Based on the analysis of tweets crawled from users, we demonstrate the potential of detecting clinical depression symptoms which emulate the PHQ-9 questionnaire clinicians use today.

Over the course of this dissertation, we examine and exploit multi-modal big (social) data to discern depressive behaviors using a wide variety of features including individual-level demographics. By developing a multi-modal framework and employing statistical techniques to fuse heterogeneous sets of features obtained through the processing of visual, textual, and user interaction data, we significantly enhance the current state-of-the-art approaches for identifying depressed individuals on social media as well as facilitate demographic inferences from social media. Besides providing insights into the relationship between demographics and mental health, our research assists in the design of a new breed of demographic-aware health interventions.

Altogether, these research topics, resulted in a framework, that when executed, will assist in identifying community-level risk and protective factors associated with the diagnosis and treatment of depression that could be an efficient means of studying patterns of access and utilization of mental health services to inform interventions.

Table of Contents

List of Figures	viii
Acknowledgements	ix
Dedication	x
1 Introduction	1
1.1 Depression: An emerging societal burden	1
1.2 Harnessing Internet and social media for public and mental health problems	2
1.3 Outline	3
2 Semi-Supervised Approach to Monitoring Clinical Depressive Symptoms in Social Media	6
2.1 Overview	6
2.2 Contributions	7
3 Studying Depressive Behavior through Visual Imagery	14
3.1 Overview	14
3.2 Contributions	16
4 Demographic Information Inference with Social Media	19
4.1 Overview	19
4.2 Contribution	21
5 Characterizing Linguistic Patterns in Depressive-indicative Language	23
5.1 Overview	23

5.2	Contribution	24
6	Conclusion	29
6.1	Summary	29
6.2	Future Work	31
	Bibliography	32
A	Contributions	40

List of Figures

2.1	Sample profile descriptions of self-declared depressed individuals	9
3.1	Sample profile descriptions of self-declared depressed individuals	15
4.1	The age distribution for depressed and control users in ground-truth dataset	20
4.2	Gender and Depressive Behavior Association (Chi-square test: color-code: (blue:association), (red: repulsion), size: amount of each cell’s contribution) .	21
5.1	Word usage difference of likely vulnerable individuals versus random profiles	24
5.2	Characterizing Linguistic Patterns in two aspects: Depressive-behavior and Age Distribution	27

Acknowledgments

I would like to express my gratitude to my advisor Prof. Pascal Hitzler for the continuous support of my research, for his patience and immense knowledge. I could not have imagined having a better mentor for this long journey.

Besides my advisor, I would like to thank all my colleagues and friends for the insightful discussions, for the sleepless nights we were working together before deadlines. Also, to the several mentors I have had whom pushed my boundaries in the new domains while doing internship under their supervision.

I express my gratitude to all the friends and collaborators, especially at Data Semantics (DaSe) Lab, my colleagues and mentors at Bosch Research, Weill Cornell Medicine, National Institute of health (National Library of Medicine), and Information Science Institute (ISI-USC).

Besides, I would like to thank Prof. Krishnaprasad Thirunarayan for his encouragement, insightful comments, and challenging questions.

Last but not least, a special thanks to my wife, for helping me in every step to see this through to the end. To my parents, who supported me spiritually.

Research reported in this publication was supported in part by NIMH of the National Institutes of Health (NIH) under award number R01MH105384-01A1.

Dedication

To my grandmother, who always wanted to call me 'doctor.'

Chapter 1

Introduction

1.1 Depression: An emerging societal burden

Depression is a highly prevalent public health concern and a major cause of disability worldwide. Depression affects 6.7% (i.e., about 16 million) Americans each year [42]. According to the World Mental Health Survey conducted in 17 countries, about 5% of people reported having at least one depressive episode in 2011 [34]. Untreated or undertreated depressive symptoms can lead to suicide and other chronic and risky behaviors such as drug or alcohol addiction [52]. More than 90% of people who commit suicide have a pre-existing diagnosis of depression [45]. Depression is also known to negatively affect daily aspects of life such as work, school, sleeping and eating habits, and family or personal relationships [15]. Recent studies also show strong associations between mental disorders and chronic diseases such as cardiovascular disease, diabetes, asthma [28], obesity, and several adverse health behaviors like smoking [6], physical inactivity, and heavy drinking [51].

Global efforts to curb depression involve identifying depressive symptoms through survey-based methods employing online questionnaires. These approaches suffer from under-representation as well as sampling bias. Survey data also exhibit problems due to temporal gaps between the data collection and dissemination of findings.

1.2 Harnessing Internet and social media for public and mental health problems

Recent years have witnessed rapid growth in the analysis of social media for studying a wide range of health problems from detecting the influenza epidemic [13] and cardiac arrest [8] to studying mood and mental health conditions [55, 57]. The widespread adoption of social media where people voluntarily and publicly express their thoughts, moods, emotions, and feelings, and share their daily struggles with mental health has not been adequately tapped into studying mental illnesses, such as depression. Insights gleaned from social media such as Twitter can be complementary to the current survey-based methods that can assist both governmental and non-governmental organizations in policy development.

The visual and textual content shared on different social media platforms like Twitter offer new opportunities for a deeper understanding of self-expressed depression both at an individual and community-level. For instance, the news headline "Twitter Fail: Teen Sent 144 Tweets Before Committing Suicide & No One Helped" highlights the need for better tools for gleaning useful insights from user generated content on social media platforms that can assist policy designers in providing resources for individuals with depressive symptoms. Recent analyses have lead to data-driven discoveries alongside the traditional hypothesis-testing social science process [3].

This thesis describes the efforts made to investigate this methodology, and focused on the following primary tasks.

1. Leveraging the textual content of social data to reliably capture clinical depression symptoms of a user over time and build a proactive and automatic depression screening tool.
2. Employing the content of posted images (colors, aesthetic, and facial presentation) for studying depressive symptoms.
3. Exploring how the choice of profile picture show any psychological traits corresponding to a depressed online persona.

4. Introduce a novel approach for studying the profiles pictures as a reliable source to glean demographic information such as age and gender, and proposed a model for community-level management of depression.
5. Studying the underlying themes among depressed individuals generated using multi-modal content that can be used to reliably detect depression.

The outline of the particular topics researched to accomplish these is provided in the next section.

1.3 Outline

This thesis is a cumulative dissertation that details the foundational research towards leveraging social data to address our country's one of the most pressing public health crises. As mentioned in the above introduction, the methods introduced here can be employed for the measurement of depression symptoms in the populations as a complementary approach for the current survey-based methods and give better understanding of the role of social media in identifying the high risk/protective socio-ecological factors, and their degree of influence in depressive behavior.

This research can be divided into four concrete topics that incrementally build towards this goal. The remainder of this dissertation is outlined as follows:

Chapter 2 presents the first research topic which introduce a statistical model that emulates traditional observational cohort studies conducted through online questionnaires by extracting, categorizing, and timely monitoring of different depression symptoms in an unobtrusive manner by modeling user-generated content in social media as a mixture of underlying topics evolving over time. The primary contributions referenced in this section are:

- [57] Amir Hossein Yazdavar, Hussein S Al-Olimat, Monireh Ebrahimi, Goonmeet Bajaj, Tanvi Banerjee, Krishnaprasad Thirunarayan, Jyotishman Pathak, and Amit Sheth. Semi-supervised approach to monitoring clinical depressive symptoms in social media.

In *Proceedings of the 2017 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining 2017*, pages 1191–1198. ACM, 2017

A semi-supervised statistical model to extract, categorize and monitor depression symptoms for continuous temporal analysis of individual’s tweets [57]

Chapter 3 discusses the usage of the content of posted images in terms of colors, aesthetic, facial presentation, and their associations with depressive symptoms.

Besides, the underlying relationships between visual and contextual content of likely depressed profiles further investigated.

Our contributions are as follows:

- [59] Amir Hossein Yazdavar, Mohammad Saeid Mahdavejad, Goonmeet Bajaj, William Romine, Amit Sheth, Amir Hassan Monadjemi, Krishnaprasad Thirunarayan, John M Meddar, Annie Myers, Jyotishman Pathak, et al. Multimodal mental health analysis in social media. *Plos one*, 15(4):e0226248, 2020
- *Analysis of the content of posted images in terms of colors, aesthetic, facial presentation, and their associations with depressive symptoms* [59]
- *Mental health analysis with social signals* [55]

Chapter 4 focuses on how demographic information inference can assist the better understanding of behavioral changes and population health.

First, we highlight the need for an approach to circumvent the inherent limitation of social signals for demographic inference (e.g. age, gender, race) that we believe is a key for gaining deeper insights for population health studies.

Then, we present our attempt to automate detection of social media users’ demographic information employing both visual content as well as textual content of social signals as not every user are willing to share their facial identity.

- *An approach for facial landmark localization, to predict gender and age from the profile as well as the shared images.* [59]

- *Demographic inference with textual content on social media* [58]

Chapter 5 focuses on characterizing the potential linguistic patterns in vulnerable user's language. We provide an in depth analysis of textual content generated by vulnerable users textual content in terms of their language aspects such as analytical thinking, clout, authenticity, and emotional tone.

- *An approach for characterizing linguistic patterns in depressive-behavior* [59]

Chapter 6 presents concluding remarks through a brief summary that highlights the overall contributions. Besides, we provide an outlook on future work.

Chapter 2

Semi-Supervised Approach to Monitoring Clinical Depressive Symptoms in Social Media

2.1 Overview

Recently, analysis of social media posts has been very successful in studying a variety of health problems [16, 10] and more specifically, much progress has been made in studying mood and mental health through social media content. These studies can be categorized into two major groups, namely, lexicon-based [37, 26] and supervised [17, 40, 18, 50]. These studies have suggested the individual's language style, emotion, ego-network, and user engagement as discriminating features to recognize the depression-indicative posts. However, the lexicon based approaches suffer from low recall and are highly dependent on the quality of the created lexicon. On the other hand, supervised approaches require labor intensive annotation of a huge dataset. Experiencing clinical depression is more than feeling down for a few days [40]. According to PHQ-9 clinical depression symptoms should persist for a few weeks. Hence, temporal monitoring of symptoms is crucial.

Inspired by that, we develop a statistical model which emulates traditional observa-

tional cohort studies conducted through online questionnaires by extracting, categorizing, and timely monitoring of different depression symptoms in an unobtrusive manner by modeling user-generated content in social media as a mixture of underlying topics evolving over time. To our knowledge, this is the first study that incorporates temporal analysis of user-generated content on social media for capturing these tell-tale symptoms.

The Diagnostic and Statistical Manual of Mental Disorders (DSM)¹ suggests that clinical depression can be diagnosed through the presence of a set of symptoms over a fixed period of time. The PHQ-9² is a nine item depression scale, which incorporates DSM-V. It can be utilized to screen, diagnose, and measure the severity of depression.

We formulated the research hypothesis is that depressed individuals discuss their symptoms on Twitter. Symptoms of depression include decreased pleasure in most activities (S1), feeling down (S2), sleep disorder (S3), loss of energy (S4), a significant change in appetite (S5), feeling worthless (S6), concentration problems (S7), hyper/lower activity (S8), and suicidal thoughts (S9). This is a top-down definition of depressive disorder through its “symptomatology”. To validate this hypothesis, we first manually examined symptoms in a random selection of 100 user profiles in our dataset. Table 2.1 illustrates a sample of anonymized tweets and their associated symptoms in PHQ-9.

For this topic, we formulated the following research questions considering these concerns and are addressed in the next section.

Q1. *How well can textual content in social media be harnessed to reliably capture clinical depression symptoms of a user over time and build a proactive and automatic depression screening tool?*

Q2. *Are there any underlying common themes among vulnerable users?*

2.2 Contributions

This section connects the individual contributions to the above research questions.

¹<https://www.psychiatry.org/psychiatrists/practice/dsm/dsm5>

²http://www.cqaimh.org/pdf/tool_phq9.pdf

Table 2.1: Sample of depressive-indicative phrases collected from tweets

Clinical Depression Symptoms	Depressive-indicative phrases in tweets
Feeling Down	"People hate me," "I am Ugly,"
Sleep disorder	"we will never sleep," "we're fuxx dead"
	"I'm that tired," "why can't I sleep"
Lassitude	"0 energy to do anything"
	"cba with work,"
Obsessed with weight	"Must not.eat," "must.be.thin"
	"94lbs, urgh I disgust myself"
	"Obsessed with my weight,"
Feeling bad about yourself	"I feel like a failure"
	"Im a piece of shix,"
Suicidal Thought	"I don't wanna wake up"
	"all my blades are so fuxx blunt"
	"Thinking hanging myself,"
	"how much blood can bleed from a cut"

[57] Amir Hossein Yazdavar, Hussein S Al-Olimat, Monireh Ebrahimi, Goonmeet Bajaj, Tanvi Banerjee, Krishnaprasad Thirunarayan, Jyotishman Pathak, and Amit Sheth. Semi-supervised approach to monitoring clinical depressive symptoms in social media. In *Proceedings of the 2017 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining 2017*, pages 1191–1198. ACM, 2017

This contribution addresses Research Question Q1 and describes a framework that is able to assess the level of depression expressed in tweets for each user profile in the dataset by integrating lexicon-based method (top-down processing) with data-driven method (bottom-up processing). Leveraging clinical articulation of depression, we build a depression lexicon that contains common depression symptoms from the established clinical assessment questionnaire such as PHQ-9 [30]. Then, the model rank the terms and compile a list of informative lexicon terms for each user and use them as seed terms to discover latent topics (depression symptoms) discussed by the subject in his/her tweets (bottom-up processing). Finally, we developed a probabilistic topic modeling over user tweets with partial supervision (by leveraging seeded clusters), named *semi-supervised topic modeling over time* (ssToT), to monitor clinical depression symptoms. We used ssToT to derive the per user topic (depression symptoms) distribution and per topic word distribution to screen and determine a

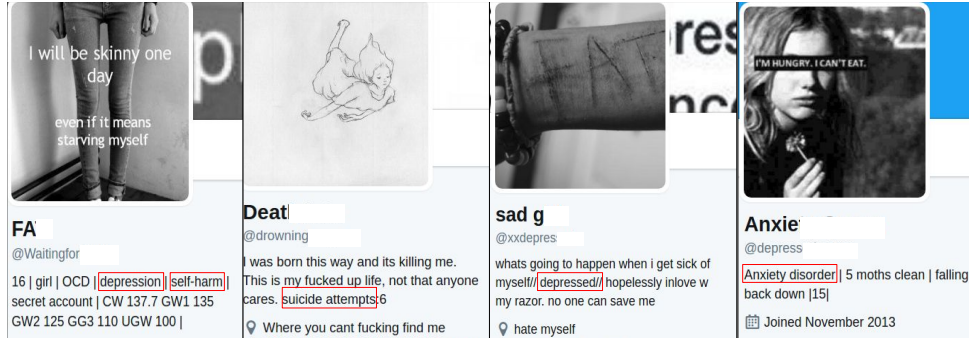


Figure 2.1: *Sample profile descriptions of self-declared depressed individuals*

trend of symptoms over time. The major contributions of this study are two-fold:

1. **First**, creating a lexicon of depression symptoms which are likely to appear in the generated content of depressed individuals.
2. **Second**, developing a semi-supervised statistical model to extract, categorize and monitor depression symptoms for continuous temporal analysis of individual’s tweets.

Dataset: We created a dataset containing 45,000 Twitter users who self-declared their depression and 2,000 “undeclared” users who were collected randomly. In particular, for collecting self-declared depressed individual’s profiles, we utilize a subset of highly informative depressive indicative terms in our lexicon and find the profiles that contain the depressive-indicative terms in their description. Figure 2.1 illustrates twitter profile descriptions of two self-declared depressed individuals.

Approach: The framework we developed aims to automatically analyze user behavior by continuously monitoring their social media content over time intervals. To this end, the approach enriches the (Latent Dirichlet Allocation) LDA model’s expressiveness by introducing a predefined set of seed terms. We divide each user’s collection of preprocessed tweets into a set of tweet buckets using a specific time interval of d days. The generative process of the proposed model for a corpus C of individual user’s tweets consisting of B buckets is shown in Algorithm-1.

In Algorithm-1, θ shows the distribution of symptoms over buckets while Φ is the distribution of words per symptom. We employ Gibbs Sampling to approximate the posterior

Algorithm 1 The generative process of ssToT

```
1: procedure ANALYZETWITTERPROFILE
2:   for each symptom (topic)  $s \in 1, 2, \dots, 9$  do
3:     Draw a distribution over terms and seed sets  $\Phi_s \sim \text{Dirichlet}(\beta)$ 
4:   end for
5:   for each bucket (document)  $b \in 1, 2, \dots, B$  do
6:     Draw a distribution over topics  $\theta_b \sim \text{Dirichlet}(\alpha)$ 
7:     for each word  $w_i \in b$  do
8:       Choose a symptom (topic)  $s_i \sim \text{multinomial}(\theta_b)$ 
9:       Choose a word  $w_i \sim \text{multinomial}(\Phi_{s_i})$ 
10:    end for
11:  end for
12: end procedure
```

distribution over the assignment of words to topics, $P(s|w)$. We then estimate Φ and θ using this posterior distribution. My strategy for discovering symptoms (topics) differs from previous methods as we incorporate prior knowledge into the inference by assigning the pre-defined seed terms into only one of the symptoms (topics). Inspired by [4], we adapt the Gibbs Sampling equation by restricting a topic s_i to a single corresponding value for each user-specific seed term or phrase. Each term w_i is assigned to the largest probability symptom associated with it in Φ . We change the probability of a symptom over a bucket to zero if the number of seed terms associated with it is less than a threshold τ . Similarly, to filter out polysemous seed terms, we aggregate the sentiment polarity of all sentences containing all seed terms over a bucket. If the aggregated polarity is positive, we assign the probability of zero to all symptoms in that bucket. Finally, we visualize the probability of each symptom over the bucket in matrix θ for further analysis and monitoring. Apart from that, if the probability of a symptom is more than a threshold τ , the symptom would be assigned to the bucket as a label. In this manner, our model can be utilized as a multi-label classifier over a time interval.

Quantitative Results Since ssToT is based on semi-supervised learning, it can be considered a clustering based approach and evaluated based on clustering evaluation measures. In addition, we are also able to employ classification accuracy to evaluate the performance of our method for symptom discovery.

Coherence Measures: Topic coherence measures score a single topic by identifying the degree of semantic similarity between high-scoring words in that topic. In this manner, we can distinguish between semantically interpretable topics and those which are artifacts of statistical inference. The state of the art for this evaluation criterion can be grouped into two major categories: intrinsic and extrinsic measures. Intrinsic measures evaluate the amount of information encoded by the topics over the original corpus used to train the topic models.

Another common intrinsic measure is UMass presented by [36], which measures the word co-occurrence in documents:

$$UMass(w_i, w_j) = \log \frac{D(w_i, w_j) + \epsilon}{D(w_i)}$$

where $D(w_i, w_j)$ counts the number of documents containing both w_i and w_j words, and $D(w_i)$ counts the ones containing w_i , over the same training corpus, and ϵ is the smoothing factor. The UMass metric computes these probabilities over the same corpus used to train the topic models.

Conversely, extrinsic evaluation metrics estimate the word co-occurrence statistics on external datasets such as Wikipedia. The UCI metric introduced by [38] utilizes the Pointwise Mutual Information (PMI) between two words,

$$UCI(w_i, w_j) = \log \frac{p(w_i, w_j) + \epsilon}{p(w_i)p(w_j)}$$

where the word probabilities are calculated by counting word co-occurrence in a sliding window over an external dataset such as Wikipedia. Recently, another topic coherence measurement has been introduced by [2] which considers context vectors for every topic's top word. For every word w , a context vector is generated using word co-occurrence counts employing context window of size $\pm n$ surrounding that word. By calculating Normalized PMI (NPMI), they showed their method has a strong correlation with human topic coherence rating. The higher the topic coherence measure score, the higher the quality of the topics. This, in turn, leads to better topic interpretability, given that our purpose is to extract meaningful and interpretable topics associable with depressive symptoms.

Table 2.2: Sample of learned topics (symptoms) by ssToT vs LDA and their associated coherency scores

Model	Label	Top Words	UMass	UCI	NPMI
ssToT	Sleep Disorder	<i>cant sleep, wanna sleep, night, nighttime, sleepy, need to sleep, hour, sleepover, bedtime, go to sleep, mess, dream, midnight crying, painful, 5:00 AM, guilt, struggle, headaches tonight, morning, coffee, duvet, hungover, bbe</i>	-1.23	-1.28	-0.01
	Eating Disorder	fat , eat, kg, weight loss, negative calories, lbs, thin, my thighs, paper thin, binge, eating disorder , abs, stomach, bulimic , hating, salad, pretend, gain, starve , mcdonalds, bones, chubby, flat, skip, wears, kcal, puffy , hippo, mfp, ugw	-1.18	0.20	0.02
	Suicidal Thoughts	self harm , cut, suicide , live, scar, blade, dead, alive, bleed, need my blade, death , hanging, deserve pain, kill me now, gun, want to die , knife, daisies, opinion, meh, razor, sharp, wrists, pictures, never wake up, wanna cut, stfu, ew	-0.66	1.13	0.09
LDA	Topic 8	<i>Sleepover, september, lost, interest, exaggerating, its my fault, ugh, sleep, skin, dish, saved, wake up, blocked, blow, ipad, touches</i>	-1.44	-2.84	-0.1
	Topic 6	<i>thigh, blood, big, beautiful, thin, smile, sleep, blood, leave, stay, worthless, fat, tear, pretending, sadness, fake, ugly, god, skin, eat, morning</i>	-3.31	-2.65	-0.08
	Topic 3	<i>Blade, ugly, fat, blood, smile, mirror, call, fit, eat, stay, beautiful, sleep, big, tear, sad, devil, god, skin, music</i>	-2.69	-3.69	-0.09

For addressing the the Research Question Q2 we conduct extensive qualitative analysis to further explore the common themes among vulnerable users.

Qualitative Results: *Discovery of common depressive symptoms*

The proposed ssToT model discovers depressive symptoms as latent topics from sliding window on buckets of timestamped tweets posted by users. We rank the top terms in each symptom $p(w|s)$ in descending order. Table 2.2 illustrates the sample of topics learned by ssToT and LDA model. The seeded words for the ssToT model are boldfaced, and words

that are judged as relevant are italicized. We observe that by constraining seed terms to a specific symptom, the discovered terms are more relevant to that category. For example, in LDA model Topic 8 contains three terms relevant to “Sleep Disorder” (S3); however, it also contains lots of irrelevant terms which makes the emphasis of this topic off-target. Although Topic 6 from LDA contains terms relevant to “Eating Disorder” (S5), it also contains some terms related to “Sleep Disorder” and “Suicidal Thoughts” (S9). Similarly, for Topic 3, it contains terms associated with both the “Eating Disorder” and “Suicidal Thoughts” categories. Therefore, the topics discovered with LDA are not interpretable for the purpose of this study.

In contrast, the topics learned from the ssToT model contain more relevant terms associated with symptom category and more interpretable topics (see Table 2.2). Additionally, the ssToT model also captures acronyms that people use in social media; for instance, in symptom 5 (Eating Disorder) “ugw” stands for “Ultimate Goal Weight” and “mfp” for “More Food Please”, or in symptom 2 (Lack of Interest) “idec” for “I Don’t Even Care”. We also observe the excessive usage of expressive interjections in language used by depressed users. Terms such as “argh” (showing frustration), “aw” (indicative of disappointment), “feh” (indicative of feeling underwhelmed), “ew” (denoting disgust), “Huh” (indicator of confusion), “phew” (showing relief) were mostly discovered in their related symptoms category.

Furthermore, we observe that there are common themes and triggers of clinical depression at the community level that they do not exist in PHQ-9. In most cases, depressed users discuss their family and friend problems and the need for their support. For example the topic {family, hugs, attention, parents, competition, daddy, mums, sigh, grandma, losing, maam, friendless, love, friend, mommy, people, boyf, gf} shows that the person is suffering from a relationship problem. Another common theme is school and academic stress {schools, college, exam, classmate, friendless, teacher, assignment}.

In summary, our proactive and automatic screening tool is able to identify clinical depressive symptoms with an accuracy of 68% and precision of 72%.

Chapter 3

Studying Depressive Behavior through Visual Imagery

3.1 Overview

According to eMarketer[33], photos accounted for 75% of content posted on Facebook worldwide, and are the most engaging type of content (87%). The ease and naturalness of expression through visual imagery can serve to glean depressive symptoms in vulnerable individuals who often seek social support through social media [48]. Thus, the choice of profile image can be a proxy for one’s online persona [31], providing a window into an individual’s mental health status. For instance, choosing a profile image with the emaciated legs of an individual with several cuts portrays negative self-view. Moreover, psychologists have argued that people use pictures to communicate messages in social media posts which represent our “Ideal Self”, or who we want to be. Indeed, we are constantly motivated to pursue behaviors that bring us closer to our Ideal Self.

The recent emergence of photo-sharing platforms such as Instagram, provides a unique opportunity to study people’s behavior through the emotions [54] with broader application in personality prediction [41] and demographic inferences. Utilizing these platforms for population-levels analysis helps to improve public health concerns [23] such as obesity [35],



Figure 3.1: *Sample profile descriptions of self-declared depressed individuals*

substance use [24], depression, and anxiety [53]. The strong associations between color sensitivity and mood has been highlighted by several studies [7]. In an earlier research, a strong correlation between specific color selection such as yellow and depressive behavior has been reported by [32]. The general findings suggest that people suffering from depression are likely to reveal their mood through their choice of colors (such as preference for darker shades) in everyday life situations [9].

Thus, the extraction of emotional state from the visual content of posted images and profile images where it can express users' emotions more vividly investigated in this research. As such, this research topic is concerned with the following questions.

Q3 . *How well does the content of posted images (colors, aesthetic, and facial presentation) reflect depressive symptoms?*

Q4 . *Does the choice of profile picture show any psychological traits corresponding to a depressed online persona?*

Table 3.1: Facial Presence Comparison in Profile/Posted images for Depressed and Control Users – * alpha = 0.05**

Face Found in	% Of Users		χ^2
	Depressed	Control	
Media	72%	81%	163.52***
Profile	4%	12%	167.2***
Not_found	8%	7%	2.55

3.2 Contributions

[59] Amir Hossein Yazdavar, Mohammad Saeid Mahdavinejad, Goonmeet Bajaj, William Romine, Amit Sheth, Amir Hassan Monadjemi, Krishnaprasad Thirunarayan, John M Meddar, Annie Myers, Jyotishman Pathak, et al. Multimodal mental health analysis in social media. *Plos one*, 15(4):e0226248, 2020

In the publication, *Multimodal Mental Health Analysis in Social Media* [59], we investigate how the visual content in posted images and profile images provide valuable psychological cues for understanding a user’s depression status. As opposed to a typical computer vision framework for object recognition that relies on thousands of predetermined low-level features, emotions reflected in facial expressions are important when assessing user’s online behavior, attributes contributing to the computational aesthetics, and sentimental quotes they may subscribe to. Besides, for capturing facial presence, we employed the model has been introduced in [61] where a multilevel convolutional coarse-to-fine network cascade developed to tackle facial landmark localization problem. We identified facial presentation and emotion from facial expression. Table 3.1 illustrates facial presentation differences in both profile and posted images (media) for depressed users and control users in U_t . For the control class, facial presence was significantly higher in both profile images and shared media (8%, 9% respectively) compared to the depressed class. In contrast with age and gender disclosure, vulnerable users were less likely to disclose their facial identity, possibly due to lack of confidence or fear of stigma.

General Image Features: The importance of interpretable computational aesthetic features for studying users’ online behavior has been highlighted before [14]. *Color*, as

a pillar of the human vision system, has a strong association with conceptual ideas like emotion [25]. We measured the normalized red, green, blue, the mean of the original colors, brightness, and contrast relative to variations of luminance. We represented images in *Hue-Saturation-Value* color space that seems intuitive for humans, and measured the mean and variance for saturation and hue. *Saturation* is defined as the difference in intensity between different light wavelengths that compose the color. Although hue is not interpretable, high saturation indicates vividness and chromatic purity, which are more appealing to the human eye [31]. *Colorfulness* is measured as a difference against gray background [46]. *Naturalness* is a measure of correspondence between images and human perception of reality [46]. In color reproduction, *naturalness* is measured from the mental recollection of the colors of familiar objects.

Additionally, there is a tendency among vulnerable users to share sentimental quotes bearing negative emotions. We performed optical character recognition (OCR) with `pythontesseract` to extract text and their sentiment [20] score. As illustrated in Table 3.2, vulnerable users tend to use less colorful (higher grayscale) profile images and shared images to convey their negative feelings, and also share images that are less natural.

In general, control users identified darker, grayer colors with negative mood, and generally preferred brighter, more vivid colors. By contrast, vulnerable users were found to prefer darker, grayer, and bluer colors. We found a strong positive correlation between self-declared depression and a tendency to perceive one’s surroundings as gray or lacking in color. With respect to the aesthetic quality of images (saturation, brightness, and hue), there is a significant difference between the two classes, with depressed users more frequently sharing images that are less appealing to the human eye.

We employed an independent samples t-test, while adopting Bonferroni Correction as a conservative approach to adjust the confidence intervals. Overall, we had 223 features, and chose Bonferroni-corrected *alpha* level of $0.05/223 = 2.24e - 4$ (***) $p < alpha$, $**p < 0.05$).

In general, the control users identified darker, grayer colors with negative moods, and generally preferred brighter, more vivid colors. In contrast, vulnerable users preferred darker, grayer colors, and bluer images. Vulnerable users shared images that are less aesthetically

Table 3.2: Statistical significance (t-statistic) of the mean of salient features for both depressed and control classes – ** alpha= 0.05, * alpha = 0.05/223**

	Feature	Depressed (μ)	Control (μ)	95 percent Conf. interval	T-stat
Image-based	Prof._colorfulness	108.05	118.85	(-15.38, -6.22)	-4.62***
	Prof._avgRGB	134.39	139.00	(2.3 6.92)	-3.92***
	Prof._naturalness	0.37	0.61	(-0.304, -0.192)	-12.72***
	Prof._hueVAR	0.0517	0.072	(-0.027, -0.008)	-4.56***
	Prof._saturationV	0.032	0.040	(-0.015, -0.003)	-3.92***
	Prof._saturationM	0.21	0.31	(-0.122, -0.078)	-8.95***
	Shared_BlueChanM	119.53	134.09	(-9.82, -19.28)	-6.04***
	Shared_GrayScaleM	0.54	0.49	(0.03, 0.068)	5.47***
	Shared_Colorfulness	106.12	122.37	(-14.98, -10.753)	-11.94***
	Shared_SaturationV	0.033	0.047	(-0.01, -0.010)	-9.26***
	Shared_SaturationM	0.198	0.289	(-0.106, -0.074)	-10.95***
	Shared_Naturalness	0.486	0.651	(-0.193, -0.136)	-16.28***
Social-based	Friends_count	610.196	1380.25	(-1023, -516)	-5.98***
	Followers_count	589.47	1340.83	(-1148.08, -354)	-3.727**
	Statuses_count	3722	7766	(-6281, -1806)	-3.55**
	Avg_favorite_cnt	0.22	0.67	(-0.781, -0.103)	-2.57**
	Avg_retweet_cnt	876.75	2720	(-2673, -1013)	-4.36***
	Favourites_count	2021	5199.67	(-5038, -1317)	-3.35**

pleasing with lower sharpness, and those that do not contain faces or contain only one face. On the other hand, control users tended to use sharper images with multiple faces. Additionally, vulnerable users shared images with more text content, often containing depressive quotes and negative sentiments. The desire to socialize and connect with others is also manifested in the visual imagery of vulnerable users. The images shared by vulnerable users tend to contain a single face (belonging to the user), rather than surrounded by friends and family. This further indicates the focus on the self, which is one of the most consistent markers of a mental disorder. This is also associated with an extensive usage of first person singular pronouns – which is another reliable marker of depression in content analysis of depressive behavior.

Chapter 4

Demographic Information Inference with Social Media

4.1 Overview

Inferring demographic information like gender and age can be crucial for stratifying our understanding of population-level epidemiology of mental health disorders. Relying on electronic health records data, previous studies have explored gender differences in depressive behavior from different angles including prevalence, age of onset, comorbidities, as well as biological and psychosocial factors. For instance, women have been diagnosed with depression twice as often as men, [43]. On the other hand, suicide rates for men are three to five times higher compared to women [5]. Although depression can affect anyone at any age, the signs and risk factors for depression vary for different age groups [12]. Depression triggers for children include domestic violence, and loss of a pet, or family member. For adolescents, depression may arise from hormonal imbalances.

Age Enabled Ground-truth Dataset: We extracted a user’s age by applying regular expression patterns to profile descriptions (such as ”17 years old, self-harm, anxiety, depression”). We compiled ”age prefixes” and ”age suffixes”, and used three age-extraction rules: 1. I am X years old, 2. Born in X, and 3. X years old, where X is a ”date” or age (e.g., 1994).

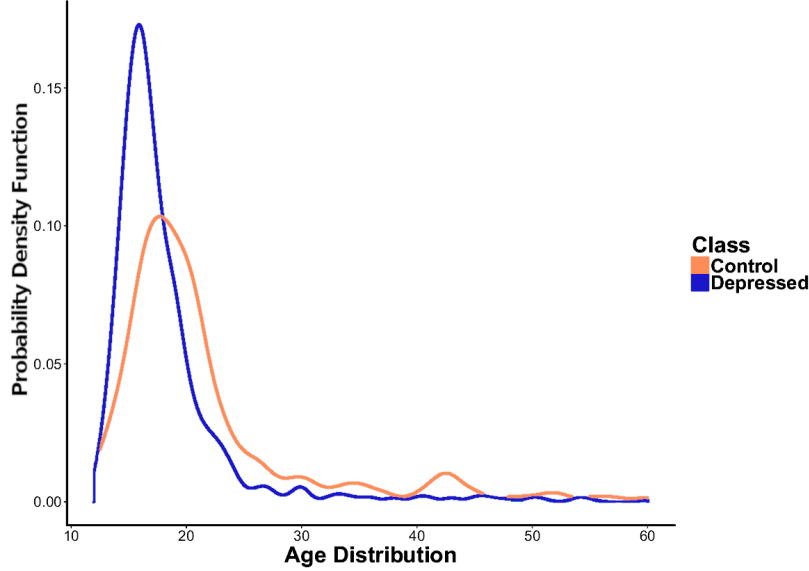


Figure 4.1: *The age distribution for depressed and control users in ground-truth dataset*

Fig 4.1 depicts the age distribution in U_a . The general trend, consistent with the results in [60], is biased toward younger individuals. Indeed, according to the Pew Research Center, 47% of Twitter users are in general 30 years old or younger [19]. The median age is 17 for the depressed class versus 19 for the control class. This suggests that the depressed-user population is younger, or depressed adolescents are more likely to disclose their age in order to connect with peers (social homophily)[1].

Gender Enabled Ground-truth Dataset: We selected a subset of 1464 users U_g from U_t who disclosed their gender in their profile description. For statistical significance, we performed a chi-square test (null hypothesis: gender and depression are two independent variables). Fig 4.2 illustrates gender association with each of the two classes. Blue circles (positive residuals, see Fig 4.2-A,D) show a positive association among corresponding row and column variables, and the red circles (negative residuals, see Fig 4.2-B,C) imply a repulsion. Our findings indicate a strong association (Chi-square: 32.75, p-value:1.04e-08) between female gender, and expression of depressive symptoms on Twitter. These observations are consistent with the current literature which have shown that more women than men are diagnosed with depression [22].

Q5 . *Are profiles pictures reliable enough to represent demographic information such as age*

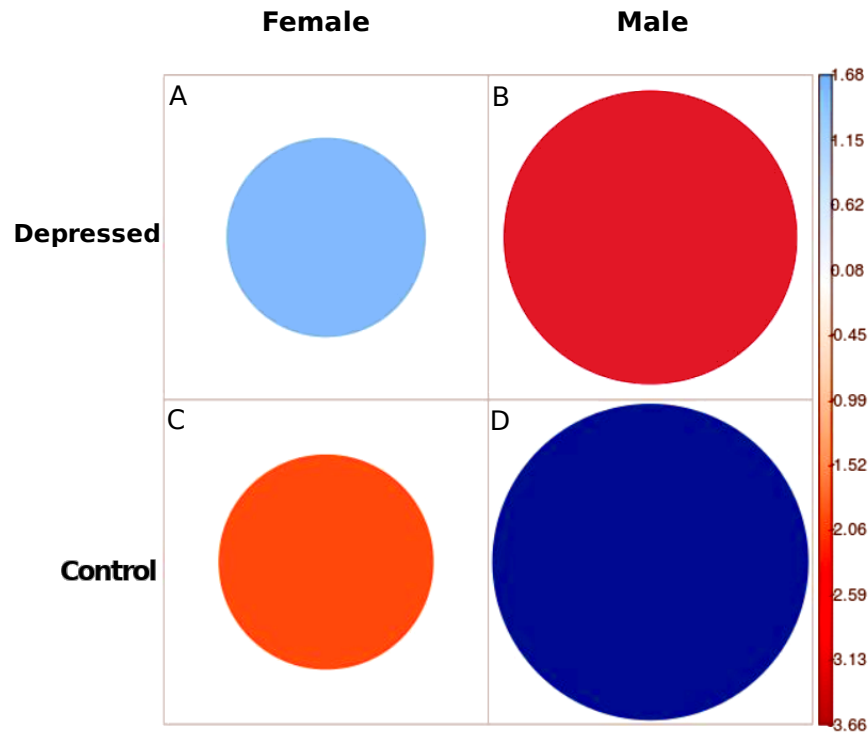


Figure 4.2: Gender and Depressive Behavior Association (Chi-square test: color-code: (blue:association), (red: repulsion), size: amount of each cell’s contribution)

and gender, and can they be used for community-level management of depression?

Q6 . Is there any other clue in the written content of the users that can be used for predicting age and gender information?

4.2 Contribution

[59] Amir Hossein Yazdavar, Mohammad Saeid Mahdavinejad, Goonmeet Bajaj, William Romine, Amit Sheth, Amir Hassan Monadjemi, Krishnaprasad Thirunarayan, John M Meddar, Annie Myers, Jyotishman Pathak, et al. Multimodal mental health analysis in social media. *Plos one*, 15(4):e0226248, 2020

In the publication, *Multimodal Mental Health Analysis is Social Media* [59], we investigate both the visual and textual content for predicting age and gender.

Prediction with Textual Content: We employed [47]’s weighted lexicon of terms that

uses the dataset of 75,394 Facebook users who shared their status, age, and gender. The predictive power of this lexica was evaluated on Twitter, and Facebook, showing promising results [47]. Utilizing these two weighted lexicon of terms, we predict the demographic information (age or gender) of $user_i$ (denoted by $Demo_i$) using the following equation:

$$Demo_i = \sum_{term \in lex} Weight_{lex}(term) * \frac{Freq(term, doc)_i}{WC(doc)_i}$$

where $Weight_{lex}(term)$ is the lexicon weight of the term, and $Freq(term, doc)_i$ represents the frequency of the term in the user generated doc_i , and $WC(doc)_i$ measures total word count in $(doc)_i$.

Prediction with Visual Imagery: Inspired by [61]’s approach for facial landmark localization, we used their pre-trained CNN consisting of convolutional layers, including unshared and fully-connected layers, to predict gender and age from both the *profile* and *shared images*.

Demographic Prediction Analysis We investigated the benefits and drawbacks of each data modality for demographic information prediction. For gender prediction, on average, the profile image-based predictor provided a more accurate prediction for both the depressed and the control class (0.92 and 0.90), compared to the content-based predictor (0.82). For age prediction, the textual content-based predictor (on average 0.60) outperformed both of the visual-based predictors (on average profile: 0.51, Media: 0.53). However, not every user provided facial identity on his or her account. We studied facial presentation for each age group to examine any association between age group, facial presentation, and depressive behavior. Less than 3% of vulnerable users between 11-19 years revealed their facial identity. Although the content-based gender predictor was not as accurate as the image-based predictor, it is adequate for population-level analysis.

Chapter 5

Characterizing Linguistic Patterns in Depressive-indicative Language

5.1 Overview

Language biases in social media posts can be a good representative of emotional state. Fig 5.1 illustrates the word clouds that distinguish the word usage of likely-depressed and non-depressed profiles. It is clear that depressed users often care more about their appearance. This is indicative by their usage of terms such as “pretty” and “beautiful.” They also have a tendency to talk about their family and relations using words such as *family, hugs, parents, daddy, mums, sigh, grandma, maam, friendless, love, friend, mommy, people ,boyf, and gf.* In contrast, the control users usually talk about daily events and news such as “hurricane” and “Trump”. Such differences in word usage highlight the fact that user generated words can be distinguishable features for detecting depressed user profiles.

Q7 . *Is there any language cues that characterize vulnerable users?*



Figure 5.1: Word usage difference of likely vulnerable individuals versus random profiles

5.2 Contribution

[59] Amir Hossein Yazdavar, Mohammad Saeid Mahdavinjad, Goonmeet Bajaj, William Romine, Amit Sheth, Amir Hassan Monadjemi, Krishnaprasad Thirunarayan, John M Meddar, Annie Myers, Jyotishman Pathak, et al. Multimodal mental health analysis in social media. *Plos one*, 15(4):e0226248, 2020

In this publication, *Multimodal Mental Health Analysis in Social Media* [49], we provided an in depth analysis of textual content generated by vulnerable users textual content in terms of analytical thinking, clout, authenticity, and emotional tone.

Thinking Style: The words we use to communicate can reveal our style of thinking. There are two common approaches for extracting an individual’s thinking style. First, measuring one’s natural way of trying to understand, analyze, and organize complex events has a strong association with analytical, formal, and logical thinking. Linguistic Inquiry and Word Count (LIWC) relates higher analytic thinking to more formal and logical reasoning, whereas a lower value indicates a focus on narratives. Second, cognitive processing, which measures problem solving in the mind, is captured through words such as ”think,” ”believe,” ”realize,” and ”know” and demonstrates ”certainty” in communication. High values for analytical thinking implies clarity of thought.

Critical thinking: ability is related to education [29], and is impacted by different stages of cognitive development at different ages [29]. It has been shown that older people communicate with greater cognitive complexity while comprehending nuances and subtle

differences [29]. All of these findings corroborate with our results (Table 5.1.)

We observed notable differences in raw intelligence and the ability to think analytically in depressed and control users among different age groups (see Figure 5.2- A, F and Table 5.1). Overall, vulnerable younger users do not think as logically based on their relative analytical score and cognitive processing ability. We can also observe that the differences between age groups above 35 tend to become smaller [21].

Authenticity: Authenticity measures the degree of honesty. Authenticity is often assessed by measuring present tense verbs, first person singular pronouns (e.g., I, me, my), and by examining the linguistic manifestations of false stories [39]. There is a decreasing trend in authenticity with age (see Fig 5.2-B.) Authenticity for depressed adolescents is strikingly higher than their control peers, and decreases with age (Fig 5.2-B.)

Clout: People with high clout speak more confidently and with certainty, employing more social words with fewer negations (e.g., no, not) and swear words. In general, mid-life is relatively stable w.r.t. relationships and work. We see the same pattern in our data (see Fig 5.2-C and Table 5.1). Unsurprisingly, lack of confidence (the 6th PHQ-9 [30] symptom) is a distinguishable characteristic of vulnerable users, leading to their lower clout scores, especially among depressed users younger than 34 years old.

Self-references: First person singular words often indicate interpersonal involvement, and their high usage is associated with negative affective states such as nervousness and depression [44]. Consistent with prior studies, the frequency of first person singular words for depressed users is significantly higher compared to that of the control class. Similarly to [44], adolescents tend to use more first-person (e.g. I), and second person singular (e.g. you) pronouns (Fig 5.2-G). The impact of the above phenomenon is reflected in significantly higher frequency of self-references for depressed adolescents. As with the control class, a downtrend suggests that as depressed individuals age, they make more distinctions and psychologically distance themselves from their topics.

Informal Language Markers; Swear, Netspeak: Swear lexicon includes terms such as “fu**”, “dam**”, and “shi*”. Several studies have highlighted that the use of profanity by young adults has significantly increased over the last decade [27]. We observed the same

pattern in both the depressed and the control classes (Table 5.1), with a higher rate for depressed users [17]. Psychologists have also shown that swearing may indicate that an individual is not a fragmented member of a society [27]. Depressed adolescents who show a higher rate of interpersonal involvement and relationships, have a higher rate of cursing (Fig 5.2-E). Also, Netspeak lexicon measures the frequency of terms such as ‘lol’ and ‘thx’. Although the rate is higher for the depressed class, we did not find any pattern concerning adult development.

Sexual, Body: The sexual lexicon contains terms like ”horny”, ”love”, and ”incest”, and body terms like ”ache”, ”heart”, and ”cough”. Both start with a higher rate for depressed users and decreases gradually as they age , possibly due to changes in sexual desire with age [11] (Fig 5.2-H,I and Table 5.1.)

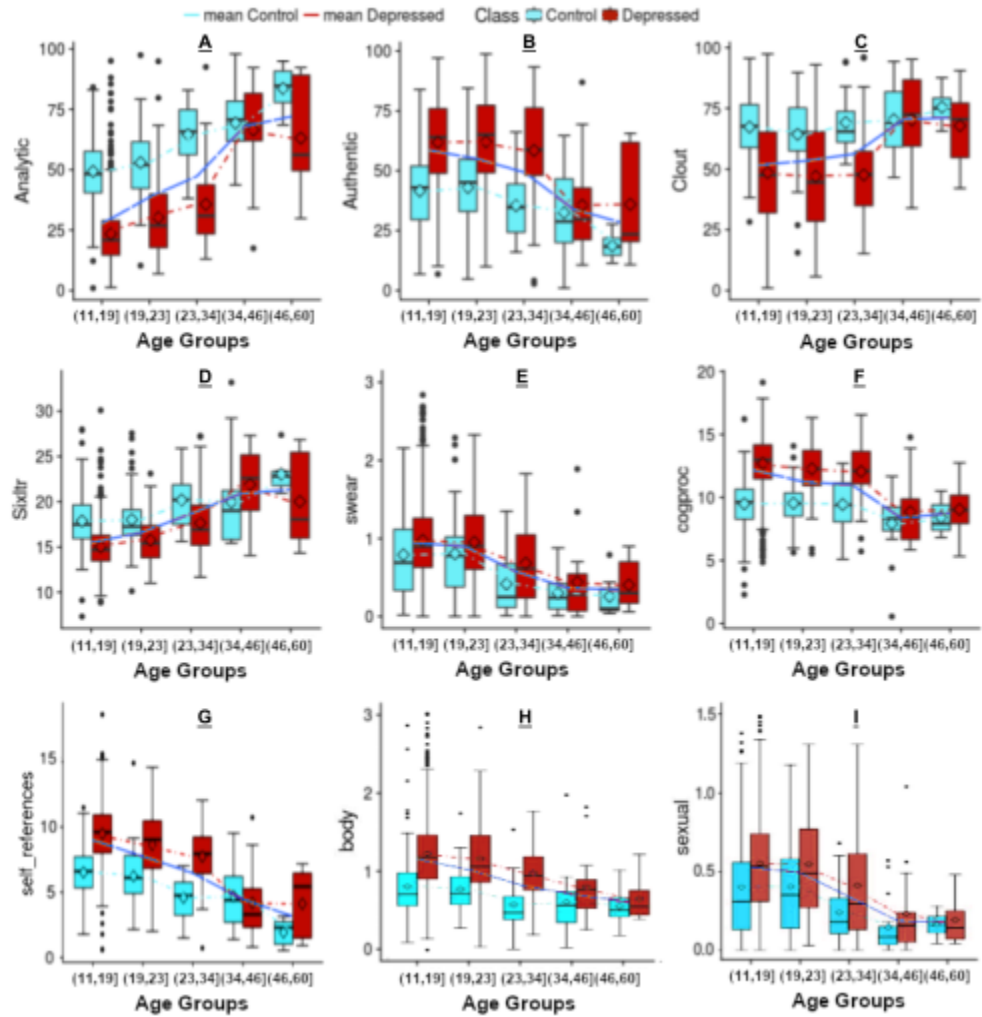


Figure 5.2: *Characterizing Linguistic Patterns in two aspects: Depressive-behavior and Age Distribution*

Table 5.1: Statistical Significance Test of Linguistic Patterns/Visual Attributes for Different Age Groups with one-way ANOVA, *** alpha = 0.001, ** alpha = 0.01

	Feature	Mean (SD)					F-value
		[11,19)	[19,23)	[23,34)	[34,46)	[46,60)	
Text-based	Analytic	27.62 (16.62)	38.61 (19.16)	47.28 (20.69)	67.88 (18.51)	72.05 (20.79)	84***
	Authentic	58.54 (19.54)	55.04 (20.04)	49.21 (22.05)	33.99 (19.73)	28.39 (19.04)	22***
	Clout	51.6 (21.35)	53.43 (21.26)	56.27 (19.81)	70.28 (17.46)	71.21 (13.50)	9***
	Dic	85.04 (6.06)	82.63 (6.21)	80.48 (6.56)	75.87 (6.91)	74.09 (5.95)	37***
	Article	3.52 (0.78)	3.92 (0.73)	4.00 (0.80)	4.52 (1.38)	5.13 (1.00)	35***
	Sixltr	15.48 (2.84)	16.58 (3.07)	18.65 (3.71)	20.88 (4.74)	21.33 (4.11)	52***
	Cogn. words	12.17 (2.53)	11.24 (2.38)	10.99 (2.55)	8.36 (2.63)	8.75 (1.96)	28***
	Self-ref	14.13 (2.35)	12.45 (2.56)	10.96 (2.60)	9.05 (3.69)	7.55 (3.38)	85***
	Swear	0.96 (0.59)	0.89 (0.53)	0.57 (0.48)	0.36 (0.41)	0.33 (0.30)	18***
	Money	0.27 (0.40)	0.38 (0.19)	0.45 (0.25)	0.52 (0.22)	0.78 (0.37)	15***
	Work	0.80 (0.39)	1.09 (0.53)	1.31 (0.76)	1.67 (0.83)	2.02 (1.01)	69***
Image-based	Prof._Naturalness	37.80 (13.84)	48.05 (18.64)	52.33 (28.51)	64.33 (24.53)	68.07 (15.28)	10***
	Prof._Saturation	20.31 (1.95)	23.27 (1.99)	29.78 (1.99)	38.76 (2.14)	33.13 (1.94)	9***
	Prof._Colorfulness	106.47 (42.70)	107.95 (39.15)	111.01 (42.09)	113.97 (35.48)	123.60 (27.60)	0.89
	Shared_avgRGB	139.20 (18.12)	140.45 (16.00)	131.55 (16.32)	133.74 (22.41)	139.02 (22.30)	3**
	Prof._GrayMean	0.471 (0.19)	0.474 (0.16)	0.456 (0.21)	0.470 (0.14)	0.450 (0.11)	0.12

Chapter 6

Conclusion

6.1 Summary

In this thesis, we demonstrated the impact of social media on extraction and timely monitoring of depression symptoms. We developed a statistical model using a hybrid approach that combines a lexicon-based technique with a semi-supervised topic modeling technique to extract per user topic distribution (clinical, symptomatic of depression) and per topic word distribution (symptom indicators) by textual analysis of tweets over different time windows. Our approach complements the current questionnaire-driven diagnostic tools by gleaning depression symptoms in a continuous and unobtrusive manner. Our experimental results reveal that there are significant differences in the topic preferences and word usage pattern of the self-declared depressed group from random users in our dataset which indicates the competency of our model for this task. Our model yields promising results with an accuracy of 68% and a precision of 72% for capturing depression symptoms per user over a time interval which is competitive with a fully supervised approach.

Besides, we presented an in-depth analysis of visual and contextual content of likely depressed profiles on Twitter. We employed them for demographic (age and gender) inference processes. We also developed a multi-modal framework, employing statistical techniques for fusing heterogeneous sets of features obtained by processing visual, textual, and user

interactions. Conducting an extensive set of experiments, we assessed the predictive power of our multi-modal framework while comparing it against state-of-the-art approaches for depressed user identification on Twitter. The empirical evaluation shows that our multi-modal framework is superior to them and it improved the average F1-Score by 5 percent. Effectively, visual cues gleaned from content and profile images shared on social media can further augment inferences from textual content for reliable determination of depression indicators and diagnoses.

Over the course of this dissertation, we have outlined a number of contributions that have improved the process to what we see today.

First, we created a semi-supervised statistical model to evaluate how the duration of these symptoms and their expression on Twitter align with the medical findings reported via the PHQ-9[57].

Next, we significantly enhanced the state-of-the-art model for identifying depressive behavior in social media by developing a multi-modal framework and employing statistical techniques to fuse heterogeneous sets of features obtained through the processing of visual, and textual content.

Besides providing insights into the relationship between demographics and mental health,our research assists in the design of a new breed of demographic-aware health interventions[56, 58].

These advances were combined, with inspiration from [56], to create a comprehensive framework that exploit different data modality for understanding depressive behavior in social media.

Altogether, these research topics, resulted in a framework, that when executed, will assisting identifying community-level risk and protective factors associated with the diagnosis and treatment of depression that could be an efficient means of studying patterns of access and utilization of mental health services to inform interventions.

6.2 Future Work

There is an immediate next step to address.

1. How can the models introduced in this research be leveraged to different data sources such as longitudinal electronic health record (EHR) systems, private insurance reimbursement, and claims data, to develop a robust “big data” platform for detecting clinical depressive behavior at the community level?

With respect to integrating the signals from social data along with longitudinal EHR data, we envision that employing the geo-location information of the social media users, we are able to perform the cross studies and expanding our findings analysis by gaining insights from the other sources.

All in all, leveraging social media for assisting public health research have a healthy outlook and can have a considerable impact on the state of the art, moving forward.

Bibliography

- [1] Faiyaz Al Zamal, Wendy Liu, and Derek Ruths. Homophily and latent attribute inference: Inferring latent attributes of twitter users from neighbors. *ICWSM*, 270:2012, 2012.
- [2] Nikolaos Aletras and Mark Stevenson. Evaluating topic coherence using distributional semantics. In *Proceedings of the 10th International Conference on Computational Semantics (IWCS 2013)–Long Papers*, 2013.
- [3] Nazanin Andalibi, Oliver L Haimson, Munmun De Choudhury, and Andrea Forte. Understanding social media disclosures of sexual abuse through the lenses of support seeking and anonymity. In *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems*, pages 3906–3918. ACM, 2016.
- [4] David Andrzejewski and Xiaojin Zhu. Latent dirichlet allocation with topic-in-set knowledge. In *Proceedings of the NAACL HLT 2009 Workshop on Semi-Supervised Learning for Natural Language Processing*. ACL, 2009.
- [5] Jules Angst, Alex Gamma, Markus Gastpar, J-P Lépine, Julien Mendlewicz, and Andre Tylee. Gender differences in depression. *European archives of psychiatry and clinical neuroscience*, 252(5):201–209, 2002.
- [6] Jafar Bakhshaie, Michael J Zvolensky, and Renee D Goodwin. Cigarette smoking and the onset and persistence of depression among adults in the united states: 1994–2005. *Comprehensive psychiatry*, 60:142–148, 2015.
- [7] Christina B Barrick, Dianne Taylor, and Elsa I Correa. Color sensitivity and mood disorders: biology or metaphor? *Journal of affective disorders*, 68(1):67–71, 2002.

- [8] Justin C Bosley, Nina W Zhao, Shawndra Hill, Frances S Shofer, David A Asch, Lance B Becker, and Raina M Merchant. Decoding twitter: Surveillance and trends for cardiac arrest and resuscitation communication. *Resuscitation*, 2013.
- [9] Helen R Carruthers, Julie Morris, Nicholas Tarrier, and Peter J Whorwell. The manchester color wheel: development of a novel way of identifying color choice and its validation in healthy, anxious and depressed individuals. *BMC medical research methodology*, 10(1):12, 2010.
- [10] Stevie Chancellor, Zhiyuan Lin, Erica L Goodman, Stephanie Zerwas, and Munmun De Choudhury. Quantifying and predicting mental illness severity in online communities. In *Proceedings of the 19th ACM Conference on Computer-Supported Cooperative Work & Social Computing*, 2016.
- [11] Eric Chung. Sexuality in ageing male: review of pathophysiology and treatment strategies for various male sexual dysfunctions. *Medical Sciences*, 7(10):98, 2019.
- [12] Mary N Cook, John Peterson, and Christopher Sheldon. Adolescent depression: an update and guide to clinical decision making. *Psychiatry (Edgmont)*, 6(9):17, 2009.
- [13] Aron Culotta. Lightweight methods to estimate influenza rates and alcohol sales volume from twitter messages. *Language resources and evaluation*, 2013.
- [14] Ritendra Datta, Dhiraj Joshi, Jia Li, and James Z Wang. Studying aesthetics in photographic images using a computational approach. In *European Conference on Computer Vision*, pages 288–301. Springer, 2006.
- [15] Munmun De Choudhury, Scott Counts, and Eric Horvitz. Social media as a measurement tool of depression in populations. In *Proceedings of the 5th Annual ACM Web Science Conference*, pages 47–56. ACM, 2013.
- [16] Munmun De Choudhury, Scott Counts, Eric J Horvitz, and Aaron Hoff. Characterizing and predicting postpartum depression from shared facebook data. In *Proceedings of*

- the 17th ACM conference on Computer supported cooperative work & social computing*, pages 626–638.
- [17] Munmun De Choudhury, Michael Gamon, Scott Counts, and Eric Horvitz. Predicting depression via social media. *Icwsn*, 13:1–10, 2013.
- [18] Virgile Landeiro Dos Reis and Aron Culotta. Using matched samples to estimate the effects of exercise on mental health from twitter. In *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence*, pages 182–188, 2015.
- [19] Maeve Duggan, Nicole B Ellison, Cliff Lampe, Amanda Lenhart, and Mary Madden. Demographics of key social networking platforms. *Pew Research Center*, 9, 2015.
- [20] Monireh Ebrahimi, Amir Hossein Yazdavar, and Amit Sheth. Challenges of sentiment analysis for dynamic events. *IEEE Intelligent Systems*, 32(5):70–75, 2017.
- [21] Penelope Eckert. Age as a sociolinguistic variable. *The handbook of sociolinguistics*, pages 151–167, 2017.
- [22] Earl S Ford, Wayne H Giles, and William H Dietz. Prevalence of the metabolic syndrome among us adults: findings from the third national health and nutrition examination survey. *Jama*, 287(3):356–359, 2002.
- [23] Venkata Rama Kiran Garimella, Abdulrahman Alfayad, and Ingmar Weber. Social media image analysis for public health. In *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems*, pages 5543–5547. ACM, 2016.
- [24] Saeed Hassanpour, Naofumi Tomita, Timothy DeLise, Benjamin Crosier, and Lisa A Marsch. Identifying substance use risk based on deep neural networks and instagram social media data. *Neuropsychopharmacology*, 44(3):487, 2019.
- [25] Kai-Qi Huang, Qiao Wang, and Zhen-Yang Wu. Natural color image enhancement and evaluation algorithm based on human visual system. *Computer Vision and Image Understanding*, 103(1):52–63, 2006.

- [26] Christian Karmen, Robert C Hsiung, and Thomas Wetter. Screening internet forum participants for depression symptoms by assembling and enhancing multiple nlp methods. *Computer methods and programs in biomedicine*, 2015.
- [27] Barbara K Kaye and Barry S Sapolsky. Watch your mouth! an analysis of profanity uttered by children on prime-time television. *Mass Communication & Society*, 7(4):429–452, 2004.
- [28] Tay Chee Kiang, Yü Anthony, Chan Kwok Wai Adrian, Lapperre Therese Sophie, Koh Mariko Siyue, et al. Anxiety, depression and hyperventilation symptoms in treatment-resistant severe asthma. *Clinical and Translational Allergy*, 5(2):P7, 2015.
- [29] Jean Kintgen-Andrews. Critical thinking and nursing education: Perplexities and insights. *Journal of Nursing Education*, 30(4):152–157, 1991.
- [30] Kurt Kroenke, Robert L Spitzer, and Janet BW Williams. The phq-9. *Journal of general internal medicine*.
- [31] Leqi Liu, Daniel Preotiuc-Pietro, Zahra Riahi Samani, Mohsen Ebrahimi Moghaddam, and Lyle H Ungar. Analyzing personality through social media profile picture choice. In *ICWSM*, pages 211–220, 2016.
- [32] Max Lüscher. *The Lüscher color test*. Simon and Schuster, 1990.
- [33] PATRICIA MARANGA. Social photos generate more engagement: New research, 2014.
- [34] Marina Marcus, M Taghi Yasamy, Mark van Ommeren, Dan Chisholm, Shekhar Saxena, et al. Depression: A global public health concern. *WHO Department of Mental Health and Substance Abuse*, 1:6–8, 2012.
- [35] Yelena Mejova, Hamed Haddadi, Anastasios Noulas, and Ingmar Weber. # foodporn: Obesity patterns in culinary interactions. In *Proceedings of the 5th international conference on digital health 2015*, pages 51–58. ACM, 2015.

- [36] David Mimno, Hanna M Wallach, Edmund Talley, and Andrew McCallum. Optimizing semantic coherence in topic models. In *Empirical Methods in NLP*. ACL, 2011.
- [37] Yair Neuman, Yohai Cohen, Dan Assaf, and Gabbi Kedma. Proactive screening for depression through metaphorical and automatic text analysis. *Artificial intelligence in medicine*, 56(1):19–25, 2012.
- [38] David Newman, Jey Han Lau, Karl Grieser, and Timothy Baldwin. Automatic evaluation of topic coherence. In *Human Language Technologies: The 2010 Conference of the North American Chapter of the ACL*, 2010.
- [39] Matthew L Newman, James W Pennebaker, Diane S Berry, and Jane M Richards. Lying words: Predicting deception from linguistic styles. *Personality and social psychology bulletin*, 29(5):665–675, 2003.
- [40] Thin Nguyen, Dinh Phung, Bo Dao, Svetha Venkatesh, and Michael Berk. Affective and content analysis of online depression communities. *IEEE Transactions on Affective Computing*, 5(3):217–226, 2014.
- [41] Jie Nie, Zhiqiang Wei, Zhen Li, Yan Yan, and Lei Huang. Understanding personality of portrait by social embedding visual features. *Multimedia Tools and Applications*, 78(1):727–746, 2019.
- [42] NIMH. How psychotherapy and other treatments can help people recover, 2014.
- [43] Susan Nolen-Hoeksema. Sex differences in unipolar depression: evidence and theory. *Psychological bulletin*, 101(2):259, 1987.
- [44] James W Pennebaker and Lori D Stone. Words of wisdom: Language use over the life span. *Journal of personality and social psychology*, 85(2):291, 2003.
- [45] M David Rudd, Alan L Berman, Thomas E Joiner Jr, Matthew K Nock, Morton M Silverman, Michael Mandrusiak, Kimberly Van Orden, and Tracy Witte. Warning signs

- for suicide: Theory, research, and clinical applications. *Suicide and Life-Threatening Behavior*, 36(3):255–262, 2006.
- [46] Jose San Pedro and Stefan Siersdorfer. Ranking and classifying attractiveness of photos in folksonomies. In *Proceedings of the 18th international conference on World wide web*, pages 771–780. ACM, 2009.
- [47] Maarten Sap, Gregory Park, Johannes Eichstaedt, Margaret Kern, David Stillwell, Michal Kosinski, Lyle Ungar, and Hansen Andrew Schwartz. Developing age and gender predictive lexica over social media. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1146–1151, 2014.
- [48] Elizabeth M Seabrook, Margaret L Kern, and Nikki S Rickard. Social networking sites, depression, and anxiety: a systematic review. *JMIR mental health*, 3(4), 2016.
- [49] Cogan Shimizu. Towards a comprehensive modular ontology IDE and tool suite. In Sabrina Kirrane and Lalana Kagal, editors, *Proceedings of the Doctoral Consortium at ISWC 2018 co-located with 17th International Semantic Web Conference (ISWC 2018), Monterey, USA, October 8th–12th, 2018*, volume 2181 of *CEUR Workshop Proceedings*, pages 65–72. CEUR-WS.org, 2018.
- [50] Hong-Han Shuai, Chih-Ya Shen, De-Nian Yang, Yi-Feng Lan, Wang-Chien Lee, Philip S Yu, and Ming-Syan Chen. Mining online social data for detecting social network mental disorders. In *Proceedings of the 25th International Conference on World Wide Web*, pages 275–285, 2016.
- [51] Tara W Strine, Ali H Mokdad, Lina S Balluz, Olinda Gonzalez, Raquel Crider, Joyce T Berry, and Kurt Kroenke. Depression and anxiety in the united states: findings from the 2006 behavioral risk factor surveillance system. *Psychiatric Services*, 59(12):1383–1390, 2008.

- [52] Lynn E Sullivan, David A Fiellin, and Patrick G O'Connor. The prevalence and impact of alcohol problems in major depression: a systematic review. *The American journal of medicine*, 118(4):330–341, 2005.
- [53] John Torous, Mark E Larsen, Colin Depp, Theodore D Cosco, Ian Barnett, Matthew K Nock, and Joe Firth. Smartphones, sensors, and machine learning to advance real-time prediction and interventions for suicide prevention: a review of current progress and next steps. *Current psychiatry reports*, 20(7):51, 2018.
- [54] Yilin Wang, Suhang Wang, Jiliang Tang, Huan Liu, and Baoxin Li. Unsupervised sentiment analysis for social media images. In *IJCAI*, pages 2378–2379, 2015.
- [55] A. H. Yazdavar, M. S. Mahdavinejad, G. Bajaj, K. Thirunarayan, J. Pathak, and A. Sheth. Mental health analysis via social media data. In *2018 IEEE International Conference on Healthcare Informatics (ICHI)*, pages 459–460, June 2018.
- [56] Amir Hossein Yazdavar, Hussein S Al-Olimat, Tanvi Banerjee, Krishnaprasad Thirunarayan, and Amit P Sheth. Analyzing clinical depressive symptoms in twitter. 2016.
- [57] Amir Hossein Yazdavar, Hussein S Al-Olimat, Monireh Ebrahimi, Goonmeet Bajaj, Tanvi Banerjee, Krishnaprasad Thirunarayan, Jyotishman Pathak, and Amit Sheth. Semi-supervised approach to monitoring clinical depressive symptoms in social media. In *Proceedings of the 2017 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining 2017*, pages 1191–1198. ACM, 2017.
- [58] Amir Hossein Yazdavar, Mohammad Saeid Mahdavinejad, Goonmeet Bajaj, William Romine, Amirhassan Monadjemi, Krishnaprasad Thirunarayan, Amit Sheth, and Jyotishman Pathak. Fusing visual, textual and connectivity clues for studying mental health. *arXiv preprint arXiv:1902.06843*, 2019.
- [59] Amir Hossein Yazdavar, Mohammad Saeid Mahdavinejad, Goonmeet Bajaj, William Romine, Amit Sheth, Amir Hassan Monadjemi, Krishnaprasad Thirunarayan, John M

- Meddar, Annie Myers, Jyotishman Pathak, et al. Multimodal mental health analysis in social media. *Plos one*, 15(4):e0226248, 2020.
- [60] Jinxue Zhang, Xia Hu, Yanchao Zhang, and Huan Liu. Your age is no secret: Inferring microbloggers' ages via content and interaction analysis. In *ICWSM*, pages 476–485, 2016.
- [61] Erjin Zhou, Haoqiang Fan, Zhimin Cao, Yuning Jiang, and Qi Yin. Extensive facial landmark localization with coarse-to-fine convolutional network cascade. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 386–391, 2013.

Appendix A

Contributions

In the following pages, all contributions made for this dissertation are listed in order of appearance in this document.

RESEARCH ARTICLE

Multimodal mental health analysis in social media

Amir Hossein Yazdavar^{1,7*}, Mohammad Saeid Mahdavejad^{4,6}, Goonmeet Bajaj², William Romine⁵, Amit Sheth³, Amir Hassan Monadjemi⁴, Krishnaprasad Thirunarayan⁶, John M. Meddar⁷, Annie Myers⁷, Jyotishman Pathak⁷, Pascal Hitzler¹

1 Department of Computer Science, Kansas State University, KS, United States of America, **2** Department of Computer Science & Engineering, Ohio State University, OH, United States of America, **3** College of Engineering and Computing, University of South Carolina, SC, United States of America, **4** Department of Artificial Intelligence & Computer Engineering, University of Isfahan, Isfahan, Iran, **5** Department of Biological Sciences, Wright State University, OH, United States of America, **6** Department of Computer Science and Engineering, Wright State University, OH, United States of America, **7** Department of Health Care Policy and Research, Weill Cornell Medicine, Cornell University, New York, NY, United States of America

* yazdavar@ksu.edu



OPEN ACCESS

Citation: Yazdavar AH, Mahdavejad MS, Bajaj G, Romine W, Sheth A, Monadjemi AH, et al. (2020) Multimodal mental health analysis in social media. *PLoS ONE* 15(4): e0226248. <https://doi.org/10.1371/journal.pone.0226248>

Editor: Jichang Zhao, Beihang University, CHINA

Received: March 6, 2019

Accepted: November 21, 2019

Published: April 10, 2020

Copyright: © 2020 Yazdavar et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Data Availability Statement: Data is available upon the acceptance to researchers whom sign the data agreement form. However, we highlight that conducting this research considered sensitive as it is related to mental health. Thus, a responsible use of this data requires an IRB, and we are unable to make this data publicly available. However, we can provide a controlled access to this dataset by having the researchers submit data agreement which is reviewed by Weill Cornell Medicine. Inquiries may be sent to Sajjad Abedian, Research Informatics Business Analyst (saa3011@med.cornell.edu).

Abstract

Depression is a major public health concern in the U.S. and globally. While successful early identification and treatment can lead to many positive health and behavioral outcomes, depression, remains undiagnosed, untreated or undertreated due to several reasons, including denial of the illness as well as cultural and social stigma. With the ubiquity of social media platforms, millions of people are now sharing their online persona by expressing their thoughts, moods, emotions, and even their daily struggles with mental health on social media. Unlike traditional observational cohort studies conducted through questionnaires and self-reported surveys, we explore the reliable detection of depressive symptoms from tweets obtained, unobtrusively. Particularly, we examine and exploit multimodal big (social) data to discern depressive behaviors using a wide variety of features including individual-level demographics. By developing a multimodal framework and employing statistical techniques to fuse heterogeneous sets of features obtained through the processing of visual, textual, and user interaction data, we significantly enhance the current state-of-the-art approaches for identifying depressed individuals on Twitter (improving the average F1-Score by 5 percent) as well as facilitate demographic inferences from social media. Besides providing insights into the relationship between demographics and mental health, our research assists in the design of a new breed of demographic-aware health interventions.

Introduction

Depression is a highly prevalent public health concern and a major cause of disability worldwide. Depression affects 6.7% (i.e., about 16 million) Americans each year [1]. According to the World Mental Health Survey conducted in 17 countries, about 5% of people reported having at least one depressive episode in 2011 [2]. Untreated or undertreated depressive

Funding: This work was supported by NIH R01MH105384-01A1 (Jyotishman Pathak, PI, Amit Sheth, PI, <https://federalreporter.nih.gov/Projects/Details/?projectId=891050>). The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Competing interests: The authors have declared that no competing interests exist.

symptoms can lead to suicide and other chronic and risky behaviors such as drug or alcohol addiction [3]. More than 90% of people who commit suicide have a pre-existing diagnosis of depression [4].

Global efforts to curb depression involve identifying depressive symptoms through survey-based methods employing online questionnaires. These approaches suffer from under-representation as well as sampling bias. Survey data also exhibit problems due to temporal gaps between the data collection and dissemination of findings.

Recent years have witnessed rapid growth in the analysis of social media for studying a wide range of health problems from detecting the influenza epidemic [5] and cardiac arrest [6] to studying mood and mental health conditions [7, 8]. The widespread adoption of social media where people voluntarily and publicly express their thoughts, moods, emotions, and feelings, and share their daily struggles with mental health has not been adequately tapped into studying mental illnesses, such as depression. Insights gleaned from social media such as Twitter can be complementary to the current survey-based methods that can assist both governmental and non-governmental organizations in policy development.

The visual and textual content shared on different social media platforms like Twitter offer new opportunities for a deeper understanding of self-expressed depression both at an individual and community-level. For instance, the news headline “Twitter Fail: Teen Sent 144 Tweets Before Committing Suicide & No One Helped” highlights the need for better tools for gleaned useful insights from user generated content on social media platforms that can assist policy designers in providing resources for individuals with depressive symptoms. Recent analyses have led to data-driven discoveries alongside the traditional hypothesis-testing social science process [9]. They have suggested that language style, sentiment, users’ activities, and engagement expressed in social media posts can predict the likelihood of depression [10, 11]. These studies often use psycholinguistic analysis, supervised and unsupervised language modeling, and expressed topics of interest. However, except for a few attempts, [12–15], these investigations have seldom studied extraction of emotional state from the visual content of posted images and profile images. Visual content can express users’ emotions more vividly, and psychologists have noted that imagery is an effective medium for communicating difficult emotions.

According to eMarketer [16], photos accounted for 75% of content posted on Facebook worldwide, and are the most engaging type of content (87%). Indeed, “a picture is worth a thousand words” and now, “photos are worth a million likes.” Similarly, on Twitter, the tweets with image links get twice as much attention as those without [17], and video-linked tweets drive up engagement [18]. The ease and naturalness of expression through visual imagery can serve to glean depressive symptoms in vulnerable individuals who often seek social support through social media [19]. Further, as psychologist Carl Rogers highlights, we often pursue and promote our Ideal-Self. In this regard, the choice of profile image can be a proxy for one’s online persona [20], providing a window into an individual’s mental health status. For instance, choosing a profile image with the emaciated legs of an individual with several cuts portrays negative self-view [21]. Moreover, psychologists have argued that people use pictures to communicate messages in social media posts which represent our “Ideal Self”, or who we want to be. Indeed, we are constantly motivated to pursue behaviors that bring us closer to our Ideal Self.

Inferring demographic information like gender and age can be crucial for stratifying our understanding of population-level epidemiology of mental health disorders. Relying on electronic health records data, previous studies have explored gender differences in depressive behavior from different angles including prevalence, age of onset, comorbidities, as well as biological and psychosocial factors. For instance, women have been diagnosed with depression

twice as often as men, [22] and a national psychiatric morbidity survey in the UK has shown a higher risk of depression in women [23]. On the other hand, suicide rates for men are three to five times higher compared to women [24]. Women are more likely to socialize and express their dysphoria, while men tend to express their anger and show negative behaviors such as alcohol abuse and drug dependency [25].

Although depression can affect anyone at any age, the signs and risk factors for depression vary for different age groups [26]. Depression triggers for children include domestic violence, and loss of a pet, or family member. For adolescents, depression may arise from hormonal imbalances [27].

Late-life depression has caused the suicide rate in people aged 80 to 84 to be more than twice that of the general population [28]. Depression in the elderly population often occurs with other medical conditions that persist, which can increase the risk of death. Therefore, inferring demographic information while studying depressive behavior from passively sensed social data can shed better light on the population-level epidemiology of depression.

The recent advancements in deep neural networks, specifically for image analysis tasks, can lead to detecting demographic features such as age and gender [29]. We aim to show that by determining and integrating a heterogeneous set of features from different modalities— aesthetic features from posted images (colorfulness, hue variance, sharpness, brightness, blurriness, naturalness), choice of profile picture (for gender, age, and facial expression), screen name, language features from both textual content and profile's description (n-gram, emotion, sentiment), sociability from ego-network, and user engagement—we can identify individuals who are more likely to be depressed from a data set of 8,770 human-annotated Twitter users.

We address the following research questions: 1) How well does the content of posted images (colors, aesthetic, and facial presentation) reflect depressive symptoms? 2) Does the choice of profile picture show any psychological traits corresponding to a depressed online persona? 3) Are profiles pictures reliable enough to represent demographic information such as age and gender, and can they be used for community-level management of depression? 4) Are there any underlying themes among depressed individuals generated using multimodal content that can be used to reliably detect depression?

Our contributions include:

- Analysis of the content of posted images in terms of colors, aesthetic, facial presentation, and their associations with depressive symptoms;
- Uncovering the underlying relationships between visual and contextual content of likely depressed profiles obtained using a demographic inference process which can facilitate community-level management of depression; and
- Testing the performance of our interpretable heterogeneous feature set for predicting depressive symptoms.

1 Related work

We have divided the related work into four subsections. First, we discuss the state-of-the-art approaches for studying depressive behavior on social data. Second, we review studies that have inferred demographic information using social media data. Then, we discuss the association between color sensitivity and mental health disorders. Finally, we cover state-of-the-art studies that have used visual imagery to study individual's behavior.

1.1 Mental health analysis using social media

Several efforts have attempted to automatically detect depression from social media content utilizing machine learning, deep learning, and natural language processing approaches. From conducting a retrospective study of tweets, De Choudhury *et al.*, (2013) characterizes depression based on factors such as language, emotion, style, ego-network, and user engagement. They built a classifier to predict the likelihood of depression from a written post [30] or an individual's profile [31]. Moreover, there have been significant advances due to the shared task [32] focusing on methods for identifying depressed users on Twitter at the Computational Linguistics and Clinical Psychology Workshop (CLP 2015). A corpus of nearly 1,800 Twitter users was built for evaluation, and the best models employed topic modeling [33], Linguistic Inquiry and Word Count (LIWC) features, and other metadata [34]. More recently, a neural network architecture has been introduced [35] to combine Twitter posts into a representation of users' activities for detecting depressed users.

Another active line of research has focused on capturing warning signs of suicide and self-harm [36]. Through analysis of tweets posted by individuals attempting committing suicide, they indicate quantifiable signals of suicidal ideations. Moreover, the CLP 2016 [36] defined a shared task on detecting the severity of mental health from forum posts. All of these studies derive discriminative features to classify depression in user-generated content at message-level, individual-level, or community-level. The recent emergence of photo-sharing platforms such as Instagram has attracted researchers' attention to study individual's behavior from their visual narratives—ranging from mining their emotions [37], and happiness trend [38], to studying medical concerns [39]. Researchers have shown that people use Instagram to engage in social exchange and share their difficult experiences [13]. The role of visual imagery as a mechanism of self-disclosure by relating visual attributes to mental health disclosures on Instagram was highlighted by [14] where individual Instagram profiles were utilized to build a prediction framework for identifying markers of depression. The importance of data modality to understand user behavior on social media has been highlighted by [40]. More recently, a deep neural network sequence modeling approach that marries audio and text data modalities to analyze question-answer style interviews between an individual and an agent has been developed to study mental health [40]. Similarly, a multimodal depressive dictionary learning process was proposed to detect depressed users on Twitter [41]. They provide sparse user representations by defining a feature set consisting of social network features, user profile features, visual features, emotional features, topic-level features, and domain-specific features. Particularly, our choice to develop a multi-modal prediction framework is intended to improve upon previous work involving the use of images in multimodal depression analysis [41] and prior work on studying Instagram photos [15].

1.2 Demographic information inference on social media

Social media has been introduced as a critical channel to answer diverse research questions offering a wealth of data for public health research [42–44].

It can also assist in better understanding the relationship between behavioral changes and population health [45]. However, the lack of demographic indicators (e.g. age, gender, race) within the data is a major limitation for gaining deeper insights. Several research efforts have attempted to automate detection of social media users' demographic information as summarized below. For gender inference, several studies have analyzed users' tweets to detect gender differences reflected in linguistic patterns [46]), profile colors [47], names [48], profile images [49], social network connections [50], and user description [46]. For instance, a supervised model was developed by [51] to determine users' gender by employing features such as screen-

name, full name, profile description, and content on external resources (e.g., personal blog). Another supervised model was built to predict the user's age group by employing features including emoticons, acronyms, slang words and phrases, punctuation, capitalization, sentence length, and included links/images, along with online behaviors such as number of friends, post time, and commenting activity [52]. To attempt to infer the age of Dutch Twitter users, a model was built that utilizes the life stage of users such as secondary school student, college student, or employee [53]. Similarly, a novel model was introduced for extracting age for Twitter users by relying on profile descriptions while devising a set of rules and patterns [54]. They also parse descriptions for occupation by consulting the SOC2010 list of occupations [55] and validating it through social surveys. A novel age inference model was developed while relying on homophily interaction information and content to predict the age of Twitter users [56]. The intuition is that people within the same age group share similar content and become friends with contemporaries. Using an extensive set of experiments, they show that their model outperformed other state-of-the-art age inference models by leveraging online interaction and content information simultaneously. The limitations of textual content for predicting age and gender was highlighted by [57]. They distinguish language use based on social gender, age identity, biological sex, and chronological age by collecting crowdsourced signals from a game in which players (crowd) guess the biological sex and age of a user based only on their tweets. Their findings indicate how linguistic markers can be misleading (e.g., a heart represented as <3 can be misinterpreted as feminine when the writer is male). Estimating age and gender from facial images by training convolutional neural networks (CNN) for face recognition is another active line of research [58].

1.3 Colors sensitivity and depressive behavior

The strong associations between color sensitivity and mood has been highlighted by several studies [59]. In an earlier research, a strong correlation between specific color selection such as yellow and depressive behavior has been reported by [60]. With respect to color discrimination, findings based on a sample of 20 male patients, aged 18 between 45 years old with schizophrenia and manic-depressive psychosis, indicated that when their right hemisphere was depressed, the identification of color by saturation, shade, and color tone was impaired [61]. More recently, the association of color vision with bipolar disorder explored [62]. The general findings suggest that people suffering from depression are likely to reveal their mood through their choice of colors (such as preference for darker shades) in everyday life situations [63]. In this study, we leveraged the visual content shared on Twitter for studying such signals.

1.4 Social media and image analysis

The recent emergence of photo-sharing platforms such as Instagram, provides a unique opportunity to study people's behavior through the emotions [37] with broader application in personality prediction [64] and demographic inferences. Utilizing these platforms for population-levels analysis helps to improve public health concerns [39] such as obesity [65], substance use [66], depression, and anxiety [67].

With regards to personality prediction, early efforts have shown that bag-of-visual-words and Facebook profile images could predict users' personality [68]. Various sets of features have been obtained from the images of 11,736 Facebook users were extracted to build a computational model which has more predictive power than human raters for predicting similar personality traits [69].

2 Dataset

This study is focused on obtaining community-level insights about depression signs and depressive behavior. As such, even though we analyzed individual’s behavioral health information—which is considered sensitive—we utilized anonymized users in our datasets as per the approved Institutional Review Board (IRB) protocol. The study was approved and the informed consent process by Wright State University Institution review Board (SC#6258) 4.1.3.

Self-disclosure refers to revealing personal and intimate information about oneself to others, which can be therapeutic for psychological well-being [70]. Previous efforts highlight diverse modes of mental health self-disclosures on social media [12]. Self-disclosure clues have been extensively utilized for creating ground-truth data for numerous social media analytic studies such as predicting users’ demographics [54], and depressive behavior [8]. For instance, vulnerable individuals may employ depressive-indicative terms in their Twitter profile descriptions. Other individuals may share their age and gender, e.g., “16 year old suicidal girl”. We employed a large dataset of 45,000 Twitter users with self-reported depressive symptoms introduced initially in [8]. All information was obtained using advanced search API [71].

To seed the search, we created a lexicon of depressive symptoms consisting of 1,500 depressive-indicative terms with the help of clinical psychologists, and employed it to collect the Twitter profiles of individuals with self-declared depressive symptoms [72]. More specifically, the dataset provides the users’ profile information including screen name, profile description, follower/followee counts, profile image, and tweet content, which can express various depression-relevant characteristics, and determine whether a user indicates any depressive behavior. Three human judges from the Department of Psychology at Wright State University assisted us in creating this annotated dataset. We reported the inter-rater agreement as $K = 0.74$ based on Cohen’s Kappa statistics [8]. To create a robust gold standard dataset, we discarded the instances in which at least two (out of three) of our annotators did not agree about the depressive symptoms. Our final dataset contains 8770 users with 3981 depressed users, and 4789 control users that do not express any depressive symptoms in their Twitter data. This dataset U_t contains the metadata values of each user such as profile descriptions, followers_count, created_at, and profile_image_url. Table 1 illustrates a sample of depressive-indicative phrases that appear in tweets from likely vulnerable users.

Table 1. Sample of depressive-indicative phrases collected from tweets.

Clinical Depression Symptoms	Depressive-indicative phrases in tweets
Feeling Down	“People hate me,” “I am Ugly,” “I am depressed”
Sleep disorder	“we will never sleep,” “we’re fuxx dead” “I’m that tired,” “why can’t I sleep”
Lassitude	“0 energy to do anything” “cba with work,” “I just want to snuggle up all day in bed”
Obsessed with weight	“Must not.eat,” “must.be.thin” “94lbs, urgh I disgust myself” “Obsessed with my weight,” “I just want be skinny”
Feeling bad about yourself	“I feel like a failure” “Im a piece of shix,”
Suicidal Thought	“I just don’t want to wake up tomorrow morning” “all my blades are so fuxx blunt” “Thinking hanging myself,” “I’ve never been so sure about suicide” “how much blood can bleed from a cut into a vain”

<https://doi.org/10.1371/journal.pone.0226248.t001>

To further measure the robustness of our dataset, we conducted another experiment by obtaining additional annotation from our colleagues from the Department of Psychiatry at Weill Cornell Medical College. Using the following formula, we computed a statistically reliable sample size:

$$\text{SampleSize} = \frac{\frac{z^2 \times p(1-p)}{e^2}}{1 + \left(\frac{z^2 \times p(1-p)}{e^2 N} \right)}$$

where N is population size, Z is z-score, e denotes margin of error, and p represents standard deviation.

Specifically, we employed our dataset of 8770 (population size), and confidence interval of 95% (margin of error 5%) to obtain 400 users as a concrete sample size. We then randomly selected 400 users from the dataset of 8770 users to be evaluated by two additional human judges (from the Department of Psychiatry at Weill Cornell Medical College) by manually annotating whether users' content reflected depressive behavior or not. The average inter-rater agreement was (85% agreement, 0.77) based on Cohen's Kappa statistics, which denotes substantial agreement and implies the robustness of our dataset.

2.1 Age enabled ground-truth dataset

We extracted a user's age by applying regular expression patterns to profile descriptions (such as "17 years old, self-harm, anxiety, depression") [54]. We compiled "age prefixes" and "age suffixes", and used three age-extraction rules: 1. I am X years old, 2. Born in X, and 3. X years old, where X is a "date" or age (e.g., 1994). We selected a subset of 1061 users among U_i as gold standard dataset U_a who disclosed their age. From these 1061 users, 822 belonged to the depressed class, and 239 belonged to the control class. From the 3981 depressed users, 20.6% disclosed their age in contrast with only 4% (239/4789) among the control group, suggesting that self-disclosure of age is more prevalent among vulnerable users. Fig 1 depicts the age distribution in U_a . The general trend, consistent with the results in [56, 73], is biased toward younger individuals. Indeed, according to the Pew Research Center, 47% of Twitter users are in general 30 years old or younger [74]. Similar data collection procedures with comparable distribution have been used previously [56]. We discuss our approach to mitigate the impact of the bias in Section 3. The median age is 17 for the depressed class versus 19 for the control class. This suggests that the depressed-user population is younger, or depressed adolescents are more likely to disclose their age in order to connect with peers (social homophily) [75].

2.2 Gender enabled ground-truth dataset

We selected a subset of 1464 users U_g from U_i who disclosed their gender in their profile description. Out of 1464 users, 64% belonged to the depressed group, and the rest (36%) belonged to the control group. 23% of the likely depressed users disclosed their gender, which is considerably higher (12%) than that of the control class. Once again, gender disclosure varies among the two gender groups. For statistical significance, we performed a chi-square test (null hypothesis: gender and depression are two independent variables). Fig 2 illustrates gender association with each of the two classes. Blue circles (positive residuals, see Fig 2A and 2D) show a positive association among corresponding row and column variables, and the red circles (negative residuals, see Fig 2B and 2C) imply a repulsion. Our findings indicate a strong association (Chi-square: 32.75, p-value: 1.04e-08) between female gender, and expression of depressive symptoms on Twitter. These observations are consistent with the current literature which have shown that more women than men are diagnosed with depression [76]. In

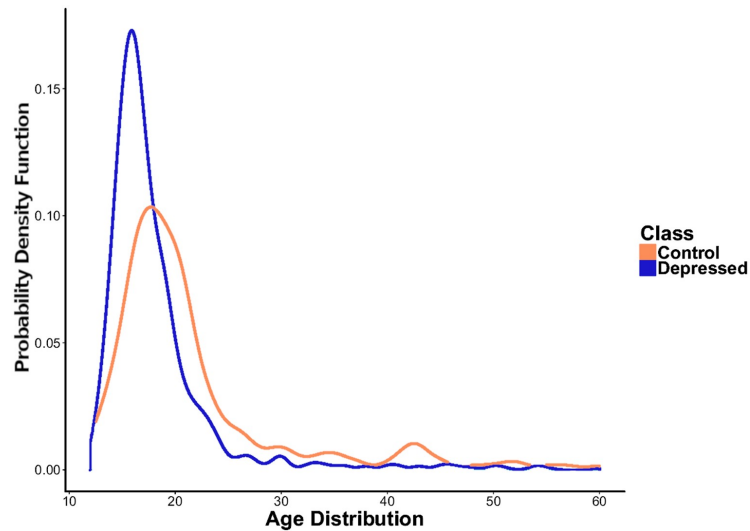


Fig 1. The age distribution for depressed and control users in ground-truth dataset.

<https://doi.org/10.1371/journal.pone.0226248.g001>

particular, the female-to-male ratio is 2:1 and 1:9 for major depressive disorder and dysthymic disorder, respectively.

3 Data modality analysis

We now provide an in-depth analysis of visual and textual content of vulnerable users.

3.1 Visual content analysis

We show that the visual content in posted images and profile images provide valuable psychological cues for understanding a user's depression status. Profile images and posted images can surface self-stigmatization [77]. As opposed to a typical computer vision framework for object recognition that relies on thousands of predetermined low-level features, emotions reflected in facial expressions are important when assessing user's online behavior, attributes contributing to the computational aesthetics, and sentimental quotes they may subscribe to.

The following sections present an in-depth analysis of visual content for both the depressed class and the control class with respect to three aspects: facial presence, facial expressions, and general image features.

3.1.1 Facial presence. For capturing facial presence, we employed the model has been introduced in [78] where a multilevel convolutional coarse-to-fine network cascade developed to tackle facial landmark localization problem. We identified facial presentation, emotion from facial expression, and demographic features from profile images and posted images [79]. Table 2 illustrates facial presentation differences in both profile and posted images (media) for depressed users and control users in U_t . For the control class, facial presence was significantly higher in both profile images and shared media (8%, 9% respectively) compared to the depressed class. In contrast with age and gender disclosure, vulnerable users were less likely to disclose their facial identity, possibly due to lack of confidence or fear of stigma.

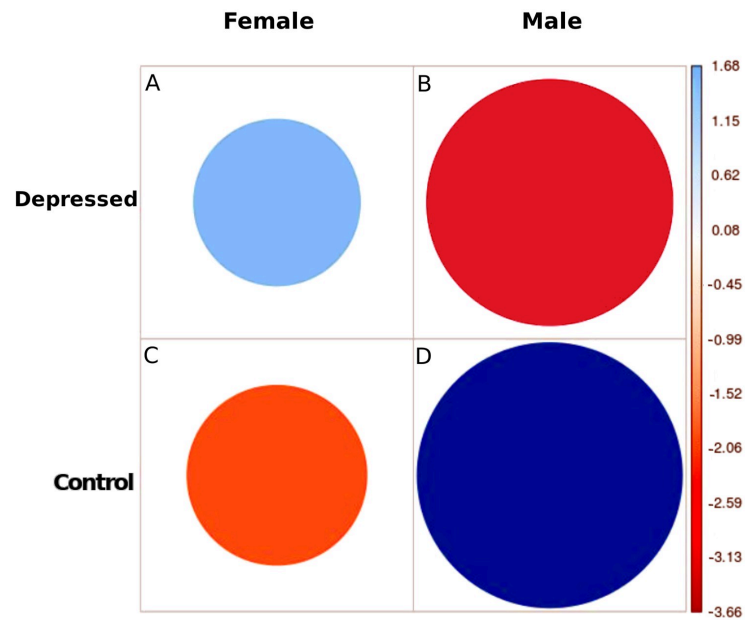


Fig 2. Gender and depressive behavior association (Chi-square test: Color-code: (blue:Association), (red: Repulsion), size: Amount of each cell's contribution).

<https://doi.org/10.1371/journal.pone.0226248.g002>

3.1.2 Facial expression. Following [20]'s approach, we adopted Ekman's model [80] of six emotions: anger, disgust, fear, joy, sadness, and surprise, and used the Face++ API [79] to automatically capture these emotions from the shared images. The positive emotions were joy and surprise, and negative emotions were anger, disgust, fear, and sadness. For each user u in U_b , we processed profile images and shared images for both the depressed and control groups with at least one face from the shared images (Table 3). For the images that contained multiple faces, we perform mean pooling over the frames to obtain the expected emotional features.

Table 2. Facial presence comparison in profile/posted images for depressed and control users—*** alpha = 0.05.

Face_Found_in	% Of Users		χ^2
	Depressed	Control	
Media	72%	81%	163.52***
Profile	4%	12%	167.2***
Not_found	8%	7%	2.55

<https://doi.org/10.1371/journal.pone.0226248.t002>

Table 3. Statistics of processed shared/profile images.

# of Processed Prof. Images		# of Processed Shared Images	
Depressed	Control	Depressed	Control
3466	4127	265785	401435

<https://doi.org/10.1371/journal.pone.0226248.t003>

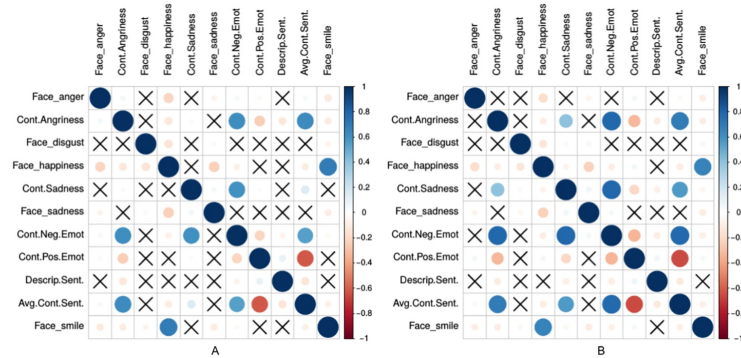


Fig 3. The Pearson correlation between the average emotions derived from facial expressions through the shared images and emotions from textual content for depressed-(a) and control users-(b). Pairs without statistically significant correlation are crossed (p-value <0.05).

<https://doi.org/10.1371/journal.pone.0226248.g003>

Fig 3 illustrates the inter-correlation of these features. Additionally, we have observed that the emotions extracted from facial expressions correlated with the emotional signals captured from textual content utilizing LIWC. This indicates that visual imagery can be utilized as a complementary channel for measuring online emotional signals.

3.1.3 General image features. The importance of interpretable computational aesthetic features for studying users’ online behavior has been highlighted by several efforts [81]. *Color*, as a pillar of the human vision system, has a strong association with conceptual ideas like emotion [82]. We measured the normalized red, green, blue, the mean of the original colors, brightness, and contrast relative to variations of luminance. We represented images in *Hue-Saturation-Value* color space that seems intuitive for humans, and measured the mean and variance for saturation and hue. *Saturation* is defined as the difference in intensity between different light wavelengths that compose the color. Although hue is not interpretable, high saturation indicates vividness and chromatic purity, which are more appealing to the human eye [20]. *Colorfulness* is measured as a difference against gray background [83]. *Naturalness* is a measure of correspondence between images and human perception of reality [83]. In color reproduction, *naturalness* is measured from the mental recollection of the colors of familiar objects. Additionally, there is a tendency among vulnerable users to share sentimental quotes bearing negative emotions. We performed optical character recognition (OCR) with python-tesseract [84] to extract text and their sentiment [85] score. As illustrated in Table 4, vulnerable users tend to use less colorful (higher grayscale) profile images and shared images to convey their negative feelings, and also share images that are less natural. In general, control users identified darker, grayer colors with negative mood, and generally preferred brighter, more vivid colors. By contrast, vulnerable users were found to prefer darker, grayer, and bluer colors. We found a strong positive correlation between self-declared depression and a tendency to perceive one’s surroundings as gray or lacking in color. With respect to the aesthetic quality of images (saturation, brightness, and hue), there is a significant difference between the two classes, with depressed users more frequently sharing images that are less appealing to the human eye.

We employed an independent samples t-test, while adopting Bonferroni Correction as a conservative approach to adjust the confidence intervals. Overall, we had 223 features, and chose Bonferroni-corrected *alpha* level of $0.05/223 = 2.24e - 4$ ($*** p < alpha, ** p < 0.05$).

Table 4. Statistical significance (t-statistic) of the mean of salient features for both depressed and control classes— alpha = 0.05, *** alpha = 0.05/223.**

	Feature	Depressed (μ)	Control (μ)	95 percent Conf. interval	T-stat
Image-based	Profile_colorfulness	108.05	118.85	(-15.38, -6.22)	-4.62***
	Profile_averageRGB	134.39	139.00	(2.3 6.92)	-3.92***
	Profile_naturalness	0.37	0.61	(-0.304, -0.192)	-12.72***
	Profile_hueVAR	0.0517	0.072	(-0.027, -0.008)	-4.56***
	Profile_saturationVAR	0.032	0.040	(-0.015, -0.003)	-3.92***
	Profile_saturationMean	0.21	0.31	(-0.122, -0.078)	-8.95***
	Shared_imageBlueChan.Mean	119.53	134.09	(-9.82, -19.28)	-6.04***
	Shared_imageGrayScaleMean	0.54	0.49	(0.03, 0.068)	5.47***
	Shared_imageColorfulness	106.12	122.37	(-14.98, -10.753)	-11.94***
	Shared_imageSaturationVAR	0.033	0.047	(-0.01, -0.010)	-9.26***
	Shared_imageSaturationMean	0.198	0.289	(-0.106, -0.074)	-10.95***
	Shared_imageNaturalness	0.486	0.651	(-0.193, -0.136)	-16.28***
Social-based	Friends_count	610.196	1380.25	(-1023, -516)	-5.98***
	Followers_count	589.47	1340.83	(-1148.08, -354)	-3.727**
	Statuses_count	3722	7766	(-6281, -1806)	-3.55**
	Avg_tweet_favorite_count	0.22	0.67	(-0.781, -0.103)	-2.57**
	Avg_tweet_retweet_count	876.75	2720	(-2673, -1013)	-4.36***
	Favourites_count	2021	5199.67	(-5038, -1317)	-3.35**

<https://doi.org/10.1371/journal.pone.0226248.t004>

In general, the control users identified darker, grayer colors with negative moods, and generally preferred brighter, more vivid colors. In contrast, vulnerable users preferred darker, grayer colors, and bluer images. Vulnerable users shared images that are less aesthetically pleasing with lower sharpness, and those that do not contain faces or contain only one face. On the other hand, control users tended to use sharper images with multiple faces. Additionally, vulnerable users shared images with more text content, often containing depressive quotes and negative sentiments.

The desire to socialize and connect with others is also manifested in the visual imagery of vulnerable users. The images shared by vulnerable users tend to contain a single face (belonging to the user), rather than surrounded by friends and family. This further indicates the focus on the self, which is one of the most consistent markers of a mental disorder. This is also associated with an extensive usage of first person singular pronouns—which is another reliable marker of depression in content analysis of depressive behavior.

3.2 Demographics inference & language cues

LIWC [86] has been used extensively for examining the latent dimensions of self-expression for analyzing personality [87], depressive behavior, demographic differences [53, 57], etc. Several studies have shown that females employ more first-person singular pronouns [88], and deictic language (context-dependent words) [89], while males tend to use more articles [90] which characterize concrete thinking, and formal, informational, affirmative words [91]. For age analysis, the salient findings show that older individuals use more future tense verbs, [88] suggesting a shift in focus while aging. They also show more positive emotions [92], employ fewer self-references (i.e. 'I', 'me'), and more first person plural pronouns [88]. Depressed users employ first person pronouns more frequently [93], and repeatedly use negative emotions and anger words. We analyzed psycholinguistic cues and language style to study the association between depressive behavior and demographics. Specifically, we adopted Levinson's adult development grouping [94] that partitions users in U_a into 5 age groups: (14,19), (19,23),

(23,34], (34,46], and (46,60]. Then, we applied LIWC for characterizing linguistic styles for each age group for users in U_n .

3.2.1 Qualitative language analysis. The recent LIWC version [86] summarizes textual content in terms of language variables such as analytical thinking, clout, authenticity, and emotional tone. It also measures other linguistic dimensions such as descriptor categories (e.g., percent of target words gleaned from the dictionary, or words longer than six letters—Sixltr), informal language markers (e.g., swear words, netspeak), and other linguistic aspects (e.g., first person singular pronouns).

Thinking Style: The words we use to communicate can reveal our style of thinking. There are two common approaches for extracting an individual's thinking style. First, measuring one's natural way of trying to understand, analyze, and organize complex events has a strong association with analytical, formal, and logical thinking. LIWC relates higher analytic thinking to more formal and logical reasoning, whereas a lower value indicates a focus on narratives. Second, cognitive processing, which measures problem solving in the mind, is captured through words such as "think," "believe," "realize," and "know" and demonstrates "certainty" in communication. High values for analytical thinking implies clarity of thought.

Critical thinking ability is related to education [95], and is impacted by different stages of cognitive development at different ages [96]. It has been shown that older people communicate with greater cognitive complexity while comprehending nuances and subtle differences [95]. All of these findings corroborate with our results (Table 5).

We observed notable differences in raw intelligence and the ability to think analytically in depressed and control users among different age groups (see Fig 4A and 4F and Table 5). Overall, vulnerable younger users do not think as logically based on their relative analytical score and cognitive processing ability. We can also observe that the differences between age groups above 35 tend to become smaller [97].

Authenticity: Authenticity measures the degree of honesty. Authenticity is often assessed by measuring present tense verbs, first person singular pronouns (e.g., I, me, my), and by examining the linguistic manifestations of false stories [98]. People who lie use fewer self-references, and fewer complex words. Psychologists often see a child's first successful lie as a mental milestone growth [99]. There is a decreasing trend in authenticity with age (see Fig 4B). Authenticity for depressed adolescents is strikingly higher than their control peers, and decreases with age (Fig 4B).

Clout: People with high clout speak more confidently and with certainty, employing more social words with fewer negations (e.g., no, not) and swear words. In general, mid-life is relatively stable w.r.t. relationships and work. A recent study has shown that age 60 is best for self-esteem [100] as people take on managerial roles at work, and maintain satisfying relationships with their spouses. We see the same pattern in our data (see Fig 4C and Table 5). Unsurprisingly, lack of confidence (the 6th PHQ-9 [101] symptom) is a distinguishable characteristic of vulnerable users, leading to their lower clout scores, especially among depressed users younger than 34 years old.

Self-references: First person singular words often indicate interpersonal involvement, and their high usage is associated with negative affective states such as nervousness and depression [92]. Consistent with prior studies, the frequency of first person singular words for depressed users is significantly higher compared to that of the control class. Similarly to [92], adolescents tend to use more first-person (e.g. I), and second person singular (e.g. you) pronouns (Fig 4G). The impact of the above phenomenon is reflected in significantly higher frequency of self-references for depressed adolescents. As with the control class, a downtrend suggests that as depressed individuals age, they make more distinctions and psychologically distance themselves from their topics.

Table 5. Statistical significance test of linguistic patterns/visual attributes for different age groups with one-way ANOVA, *** alpha = 0.001, ** alpha = 0.01.

	Feature	Mean (SD)					F-value
		[11,19]	[19,23]	[23,34]	[34,46]	[46,60]	
Text-based	Analytic	27.62 (16.62)	38.61 (19.16)	47.28 (20.69)	67.88 (18.51)	72.05 (20.79)	84***
	Authentic	58.54 (19.54)	55.04 (20.04)	49.21 (22.05)	33.99 (19.73)	28.39 (19.04)	22***
	Clout	51.6 (21.35)	53.43 (21.26)	56.27 (19.81)	70.28 (17.46)	71.21 (13.50)	9***
	Dic	85.04 (6.06)	82.63 (6.21)	80.48 (6.56)	75.87 (6.91)	74.09 (5.95)	37***
	Article	3.52 (0.78)	3.92 (0.73)	4.00 (0.80)	4.52 (1.38)	5.13 (1.00)	35***
	Sixtr	15.48 (2.84)	16.58 (3.07)	18.65 (3.71)	20.88 (4.74)	21.33 (4.11)	52***
	Cogn. words	12.17 (2.53)	11.24 (2.38)	10.99 (2.55)	8.36 (2.63)	8.75 (1.96)	28***
	Self-ref	14.13 (2.35)	12.45 (2.56)	10.96 (2.60)	9.05 (3.69)	7.55 (3.38)	85***
	Swear	0.96 (0.59)	0.89 (0.53)	0.57 (0.48)	0.36 (0.41)	0.33 (0.30)	18***
	Money	0.27 (0.40)	0.38 (0.19)	0.45 (0.25)	0.52 (0.22)	0.78 (0.37)	15***
	Work	0.80 (0.39)	1.09 (0.53)	1.31 (0.76)	1.67 (0.83)	2.02 (1.01)	69***
Image-based	Profile_Naturalness	37.80 (13.84)	48.05 (18.64)	52.33 (28.51)	64.33 (24.53)	68.07 (15.28)	10***
	Profile_SaturationMean	20.31 (1.95)	23.27 (1.99)	29.78 (1.99)	38.76 (2.14)	33.13 (1.94)	9***
	Profile_Colorfullness	106.47 (42.70)	107.95 (39.15)	111.01 (42.09)	113.97 (35.48)	123.60 (27.60)	0.89
	Shared_avgRGB	139.20 (18.12)	140.45 (16.00)	131.55 (16.32)	133.74 (22.41)	139.02 (22.30)	3**
	Profile_GrayMean	0.471 (0.19)	0.474 (0.16)	0.456 (0.21)	0.470 (0.14)	0.450 (0.11)	0.12

<https://doi.org/10.1371/journal.pone.0226248.t005>

Informal Language Markers; Swear, Netspeak: Swear lexicon includes terms such as “fu***”, “dam*”, and “shi*”. Several studies have highlighted that the use of profanity by young adults has significantly increased over the last decade [102]. We observed the same pattern in both the depressed and the control classes (Table 5), with a higher rate for depressed users [10]. Psychologists have also shown that swearing may indicate that an individual is not a fragmented member of a society [103]. Depressed adolescents who show a higher rate of interpersonal involvement and relationships, have a higher rate of cursing (Fig 4E). Also, Netspeak lexicon measures the frequency of terms such as ‘lol’ and ‘thx’. Although the rate is higher for the depressed class, we did not find any pattern concerning adult development.

Sexual, Body: The sexual lexicon contains terms like “horny”, “love”, and “incest”, and body terms like “ache”, “heart”, and “cough”. Both start with a higher rate for depressed users and decreases gradually as they age, possibly due to changes in sexual desire with age [104] (Fig 4H and 4I and Table 5).

3.2.2 Quantitative language analysis. We employed a one-way ANOVA to compare the impact of various factors, and validate our findings above. Table 5 illustrates our findings, with

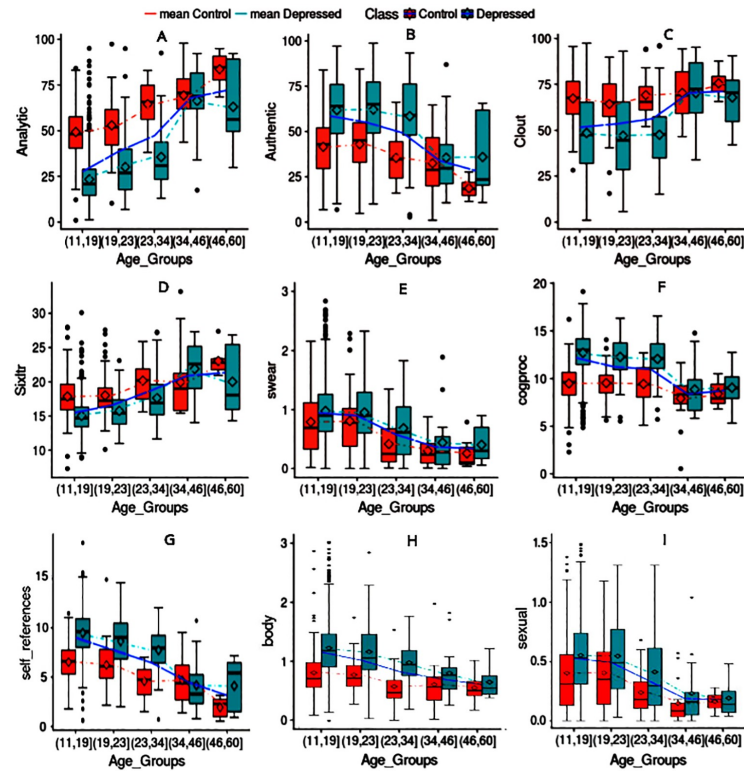


Fig 4. Characterizing linguistic patterns in two aspects: Depressive-behavior and age distribution.

<https://doi.org/10.1371/journal.pone.0226248.g004>

a degree of freedom (df) of 1055. The null hypothesis is that the sample means for each age group are similar for each of the LIWC features.

3.3 Demographic prediction

We leveraged both the visual and textual content for predicting age and gender.

3.3.1 Prediction with textual content. We employed [105]’s weighted lexicon of terms that uses the dataset of 75,394 Facebook users who shared their status, age, and gender. The predictive power of this lexica was evaluated on Twitter, and Facebook, showing promising results [105]. Utilizing these two weighted lexicon of terms, we are predicting the demographic information (age or gender) of $user_i$ (denoted by $Demo_i$) using the following equation:

$$Demo_i = \sum_{term \in lex} Weight_{lex}(term) * \frac{Freq(term, doc)_i}{WC(doc)_i}$$

where $Weight_{lex}(term)$ is the lexicon weight of the term, and $Freq(term, doc)_i$ represents the frequency of the term in the user generated doc_i , and $WC(doc)_i$ measures total word count in $(doc)_i$. As our data are biased toward younger individuals, we report age prediction

Table 6. Age Prediction performance from visual and textual content for different age group (years old).

Group	Measure	Text-based				Image-based (Profile)				Image-based (Media)			
		(11,19]	(19,23]	(23,34]	(34,46]	(11,19]	(19,23]	(23,34]	(34,46]	(11,19]	(19,23]	(23,34]	(34,46]
Depressed	Sensitivity	0.23	0.38	0.65	0.33	0.29	0.29	0.22	1.0	0.11	0.1	0.19	0.22
	Specificity	0.95	0.53	0.69	0.96	0.92	0.92	0.57	0.80	0.96	0.94	0.72	0.58
	ACC	0.59	0.46	0.67	0.65	0.47	0.46	0.40	0.900	0.50	0.49	0.46	0.40
Control	Sensitivity	0.14	0.31	0.62	0.69	0.12	0.1	0.40	0.25	0.15	0.30	0.63	0.64
	Specificity	0.98	0.63	0.61	0.90	0.90	0.95	0.53	0.75	0.98	0.62	0.60	0.91
	ACC	0.56	0.47	0.62	0.80	0.49	0.48	0.47	0.51	0.56	0.46	0.62	0.77

<https://doi.org/10.1371/journal.pone.0226248.t006>

performance for each age group, separately (Table 6). Moreover, to measure the average accuracy of this model, we built a balanced dataset (keeping the total number of users above 23—416), and then randomly sampled the same number of users from the age ranges (11,19] and (19,23]. The average accuracy of this model was 0.63 for depressed users, and 0.64 for the control class. Table 8 illustrates the performance of gender prediction for each class. The average accuracy was 0.82 on U_g ground-truth dataset.

3.3.2 Prediction with visual imagery. Inspired by [78]’s approach for facial landmark localization, we used their pre-trained CNN consisting of convolutional layers, including unshared and fully-connected layers, to predict gender and age from both the *profile* and *shared images*. We evaluated the performance of the gender and age prediction task on U_g and U_a , respectively, as shown in Table 6.

3.3.3 Demographic prediction analysis. We delved deeper into the benefits and drawbacks of each data modality for demographic information prediction. This is crucial as the differences between language cues between age groups above 35 tend to become smaller (see Fig 4A, 4B and 4C), making the prediction harder for older individuals [97]. In this case, the other data modality (e.g., visual content) played an integral role as a complementary source for age inference. For gender prediction, on average, the profile image-based predictor provided a more accurate prediction for both the depressed and the control class (0.92 and 0.90), compared to the content-based predictor (0.82). For age prediction (see Table 6), the textual content-based predictor (on average 0.60) outperformed both of the visual-based predictors (on average profile: 0.51, Media: 0.53). However, not every user provided facial identity on his or her account (see Table 2). We studied facial presentation for each age group to examine any association between age group, facial presentation, and depressive behavior (see Table 7). We can see youngsters in both the depressed and control classes are not likely to present their face in their profile image. Less than 3% of vulnerable users between 11-19 years revealed their facial identity. Although the content-based gender predictor was not as accurate as the image-based predictor, it is adequate for population-level analysis (see Table 8).

4 Multi-modal prediction framework

We used the above findings for predicting depressive behaviors. Our model exploits an early fusion [40] technique in feature space and requires modeling each user u in U_i as vector

Table 7. Facial presentation distribution for different age group (in years old) in profile and media.

	% Users Faces_Found_in_Profile					% Users Faces_Found_in_Media				
	(11,19]	(19,23]	(23,34]	(34,46]	(46,60]	(11,19]	(19,23]	(23,34]	(34,46]	(46,60]
Control	4.55	9.58	13.84	17.85	21.42	89.70	88.35	78.46	67.85	78.57
Depressed	2.71	5.88	10.52	8.33	14.28	90.21	90.58	76.31	83.33	85.71

<https://doi.org/10.1371/journal.pone.0226248.t007>

Table 8. Gender prediction performance through visual and textual content.

Face found in	Image-based Predictor						Content-based Predictor					
	Depressed			Control			Depressed			Control		
	Sens.	Spec.	ACC (95% CI)	Sens.	Spec.	ACC (95% CI)	Sens.	Spec.	ACC (95% CI)	Sens.	Spec.	ACC (95% CI)
Profile	0.90	1.0	0.92 (0.80, 0.98)	0.91	0.87	0.90 (0.81, 0.95)	0.87	0.50	0.82 (0.79, 0.85)	0.86	0.76	0.82 (0.79, 0.85)
Media	0.57	0.70	0.58 (0.54, 0.62)	0.46	0.65	0.51 (0.46, 0.55)						

<https://doi.org/10.1371/journal.pone.0226248.t008>

concatenation of individual modality features. As opposed to the computationally expensive late fusion schemes, where each modality requires a separate supervised modeling, this model reduces the learning effort and has shown promising results [106]. To develop a generalizable model that avoids overfitting, we performed feature selection using statistical tests and *all relevant* ensemble learning models. Adding feature selection tests adds randomness to the data by creating shuffled copies of all features (shadow feature), and then trains the Random Forest classifier on the extended data. Iteratively, it checks whether the actual feature has a higher Z-score than its shadow feature (See Algorithm 1 and Fig 5) [107].

Algorithm 1: Ensemble Feature Selection

```

Function Main
  for each Feature  $X_j \in X$  do
     $ShadowFeatures \leftarrow RndPerm(X_j)$ 
     $RndForrest(ShadowFeatures, X)$ ;
    Calculate  $Imp(X_j, MaxImp(ShadowFeatures))$ ;
  if  $Imp(X_j) > MaxImp(ShadowFeatures)$  then
    Generate next hypothesis, return  $X_j$ 
  Once all hypothesis generated;
  Perform Statistical Test
   $H_0 : H_i = E(H)$  vs  $H_1 : H_i \neq E(H)$   $H_i \sim N((0.5N)((\sqrt{0.25N})^2))$  //Binomial
  Distribution;
  if  $H_i \gg E(H)$  then
    Feature is important
  else
    Feature is important
  
```

Next, we adopted an ensemble learning method which integrated the predictive power of multiple learners with two main advantages; a high degree of interpretability with respect to the contributions of each feature, and a high predictive power. For prediction, we have $y'_i = \sum_{i=1}^m f_i(u_i)$ where f_i is a weak learner and y'_i denotes the final prediction.

In particular, we optimized the loss function:

$$L^{<t>} = \sum_{i=1}^n l(y_i, y_i^{<t-1>} + f_i(u_i)) + \varphi(f_i)$$

where φ incorporates L1 and L2 regularization. In each iteration, the new $f_i(u_i)$ is obtained by fitting the weak learner to the negative gradient of loss function. Particularly, by estimating the loss function with Taylor expansion:

$$L^{<t>} \sim \sum_{i=1}^n l(y_i, y_i^{<t-1>} + f_i(u_i)) + \left(\frac{\partial l(y_i, y_i^{<t-1>}}{\partial y_i^{<t-1>}}\right) f_i(u_i) + \left(\frac{\partial^2 l(y_i, y_i^{<t-1>}}{\partial y_i^{<t-1>^2}}\right) f_i(u_i)^2$$

where its first expression is constant, the second and the third expressions are first (g_i) and

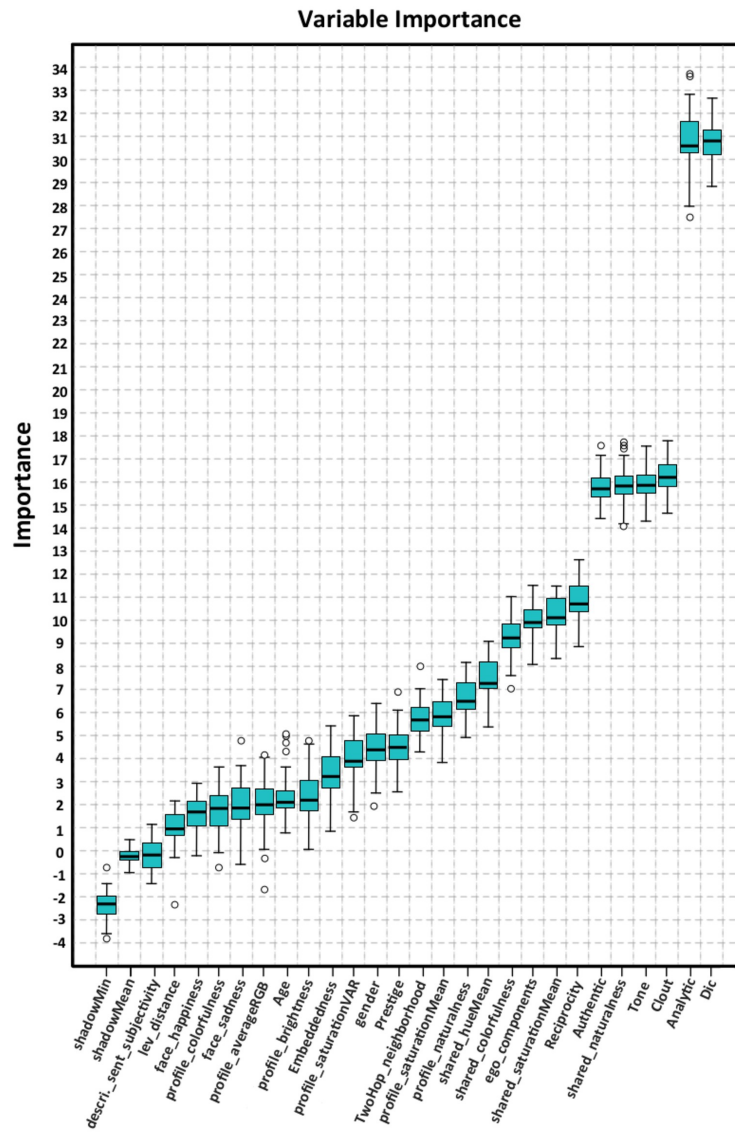


Fig 5. Ranking features obtained from different modalities with an ensemble algorithm.

<https://doi.org/10.1371/journal.pone.0226248.g005>

second order derivatives (h_i) of the loss.

$$L^{<t>} = \sum_{i=1}^n (g_i f_i(u_i) + h_i f_i(u_i)) + \varphi(f_i)$$

To explore the weak learners, assume f_i has k leaf nodes, I_j be subset of users from U_i belongs to the node j , and w_j denotes the prediction for node j . Then, for each user i belonging to I_j , $f_i(u_i) = w_j$ and $\varphi(f_i) = 1/2\lambda \sum_{j=1}^k W_j^2 + \gamma k$

$$L^{<t>} = \sum_{j=1}^k [(\sum_{i \in I_j} g_i) w_j + 1/2(\sum_{i \in I_j} h_i + \lambda) w_j^2] + \gamma k$$

Next, for each leaf node j , deriving w.r.t w_j :

$$w_j = \frac{-\sum_{i \in I_j} g_i}{\sum_{i \in I_j} h_i + \lambda}$$

and by substituting weights:

$$L^{<t>} = -1/2 \sum_{j=1}^k \frac{(\sum_{i \in I_j} g_i)^2}{\sum_{i \in I_j} h_i + \lambda} + \gamma k$$

which represents the loss of fixed weak learners with k nodes. The trees are built sequentially, such that each subsequent tree aims to reduce the errors of its predecessor trees. Although, the weak learners have a higher degree of biases, the ensemble model produces a strong learner that effectively integrates the weak learners by reducing bias and variance (the ultimate goal of supervised models) [108, 109]. Table 9 illustrates how our multimodal framework outperforms the baselines for identifying depressed users in terms of average specificity, sensitivity, F-Measure, and accuracy in a 10-fold cross-validation setting on U_i dataset. Fig 6 shows how the likelihood of being classified into the depressed class varies with each feature added to the model for a sample user in the dataset. The prediction bar (the black bar) shows that the log-odds of prediction is 0.31, that is, the likelihood of this person being a depressed user is 57% ($1 / (1 + \exp(-0.3))$). The figure also sheds light on the impact of each contributing feature. The waterfall

Table 9. Model's performance for depressed user identification in Twitter using different data modalities.

Model#	Data Source	Ref.	Year	Features					Model	Spec.	Sens.	F-1	Acc.
				N-grams	LIWC	Sentiment	Topics	Metadata					
I	Content	[110]	2016	X					NB	0.69	0.70	0.69	0.70
II		[111]	2016	X		X		User Acti.	N/A (LR)	0.73	0.74	0.73	0.74
III		[112]	2015	X	X	X		User Acti.	Log-linear	0.83	0.80	0.81	0.82
IV		[113]	2015	X	X	X	X		LR	0.84	0.83	0.84	0.84
V		[114]	2015	X	X	X	X	User Acti.	SVM	0.86	0.84	0.85	0.85
VI		N/A	N/A	X					SVM(Pre. embed.)	0.72	0.72	0.72	0.72
VII		N/A	N/A	X					SVM(Train w2vec)	0.70	0.70	0.70	0.70
VIII	Cont., Net.	[10]	2013	X	X	X			SVM, PCA	0.84	0.80	0.83	0.85
IX	Image	N/A	N/A	N/A					LR	0.68	0.67	0.67	0.68
X		N/A	N/A						SVM	0.69	0.67	0.67	0.69
XI		N/A	N/A						RF	0.72	0.70	0.69	0.71
Ours	Cont.,Image,Net.	N/A	X	X	X	X	X	X	N/A	0.87	0.92	0.90	0.90

<https://doi.org/10.1371/journal.pone.0226248.t009>

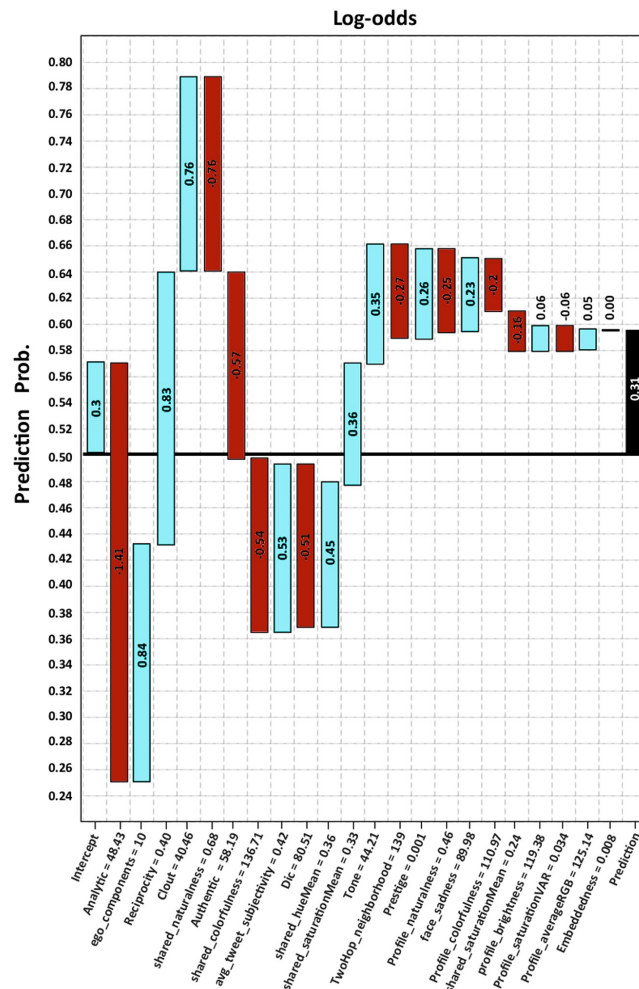


Fig 6. The explanation of the log-odds prediction of outcome (0.31) for a sample user (y-axis shows the outcome probability (depressed or control), the bar labels indicate the log-odds impact of each feature).

<https://doi.org/10.1371/journal.pone.0226248.g006>

charts represent how the probability of being depressed varies when adding each feature. For example, for our dataset, the “Analytic thinking” score measured by LIWC from the tweet content is a high value of 48.43 (Median:36.95, Mean: 40.18) and this decreases the chance of the user being classified into the depressed group by the log-odds of -1.41. This is due to the fact that depressed users have significantly lower “Analytic thinking” scores compared to the control class. Moreover, the “Clout” score of 40.46 is considered a low value (Median: 62.22, Mean: 57.17), and increases the chance of being classified as a depressed user. This is justifiable

given the clear association between low self-esteem and risk for depression. With respect to the visual features, the mean and the median of “shared colorfulness” is 112.03 and 113, respectively. The value of 136.71 would be high, and decreases the chance of being depressed by log-odds of -0.54. As mentioned earlier, depressed users preferred darker, and grayer colors. The score of 0.46 as “profile naturalness” is considered high compared to 0.36 (the mean for the depressed class) which justifies pull down of the log-odds by -0.25. Using network features, for instance, the “two hop neighborhood” for depressed users (Mean: 84) are less than that of the control users (Mean: 154), and is reflected in pulling down the log-odds by -0.27.

4.1 Baselines

To test the efficacy of our multi-modal framework for detecting depressed users, we compared it against existing content, content-network, and image-based models (based on the aforementioned general image features, facial presence, and facial expressions).

4.1.1 Content-based models. Language biases in social media posts can be a good representative of emotional state. Fig 7 illustrates the word clouds that distinguish the word usage of likely-depressed and non-depressed profiles. It is clear that depressed users often care more about their appearance. This is indicative by their usage of terms such as “pretty” and “beautiful.” They also have a tendency to talk about their family and relations using words such as *family, hugs, parents, daddy, mums, sigh, grandma, maam, friendless, love, friend, mommy, people, boyf, and gf*. In contrast, the control users usually talk about daily events and news such as “hurricane” and “Trump”. Such differences in word usage highlight the fact that user generated words can be distinguishable features for detecting depressed user profiles. See Table 9 for the comparative performance of our prediction framework against state-of-the-art methods used for predicting depressive behaviors—many of which employed the same feature sets and hyperparameter settings (see Models I-V). Several prior efforts have demonstrated that word embedding models can reliably enhance short text classification [115], Model VI by employing pre-trained word embeddings which have trained over 400 million tweets [116] while representing a user with retrieving word vectors for all the words a user used in tweets and profile description. We aggregate these word vectors through their means and feed it as input to a SVM classifier with a linear kernel. In Model VII, we employed [8]’s dataset of 45,000 self-reported depressed users and trained a Skip-gram model with negative sampling to learn word representations. We chose this model because it generates robust word embeddings even when the collection of training words are sparse [117]. We set dimensionality to 300 and a negative sampling rate to 10 sample words, which has shown promising results with medium-sized datasets [117]. Besides, we have observed that many vulnerable users chose specific



Fig 7. Word usage difference of likely vulnerable individuals versus random profiles.

<https://doi.org/10.1371/journal.pone.0226248.g007>

account names, such as “Suicidal_Thoughtxxx,” and “younganxietyyxxx,” which are good indicators of their depressive behavior. We used Levenshtein distance between depression indicative terms in [8]’s depression lexicon and the screen name to capture their degree of semantic similarity [118].

4.1.2 Image-based models. We employed the aforementioned visual content features including facial presence, aesthetic features, and facial expression for depression prediction. We use three different models: Logistic Regression (Model IX), SVM (Model X), and Random Forest (Model XI). The poor performance of image-based models suggests that relying on a unique modality would not be sufficient for building a robust model due to the complexity and abstruse nature of the prediction task.

4.1.3 Network-based models. Network-based features indicate the user’s desire to socialize and connect with others. There is a notable difference between the number of friends, followers, favorites, and status count for depressed and control users (see Table 4). For building a baseline Model VIII, we obtained egocentric network measures for each user based on the network formed using @-replies interactions among them. The egocentric social graph of a user u is an undirected graph of nodes in u ’s two-hop neighborhood in our U_a dataset, where the edge between nodes u and v implies that there has been at least one @-reply exchange. Network-based features including *Reciprocity*, *Prestige Ratio*, *Graph Density*, *Clustering Coefficient*, *Embeddedness*, *Ego components* and *Size of two-hop neighborhood* were extracted from each user’s network [10] to reliably capture user context for depression prediction.

High values for the three metrics—clustering coefficient, embeddedness, and number of ego networks—indicates that the depressed users tend to build a close network of trusted people to share their mental health issues. For both graph density and size of the two-hop neighborhood, a lower value indicates fewer interactions.

Conclusion and future work

We presented an in-depth analysis of visual and contextual content of likely depressed profiles on Twitter. We employed them for demographic (age and gender) inference processes. We developed a multimodal framework, employing statistical techniques for fusing heterogeneous sets of features obtained by processing visual, textual, and user interactions. Conducting an extensive set of experiments, we assessed the predictive power of our multimodal framework while comparing it against state-of-the-art approaches for depressed user identification on Twitter. The empirical evaluation shows that our multimodal framework is superior to them and it improved the average F1-Score by 5 percent. Effectively, visual cues gleaned from content and profile images shared on social media can further augment inferences from textual content for reliable determination of depression indicators and diagnoses. In the future, we plan to apply our approach to various data sources such as longitudinal electronic health record (EHR) systems, and private insurance reimbursement and claims data, to develop a robust “big data” platform for detecting clinical depressive behavior at the community level.

Supporting information

S1 File. The informed consent of this study approved by Wright State University Institution review Board (SC#6258).
(PDF)

Acknowledgments

Research reported in this publication was supported in part by NIMH of the National Institutes of Health (NIH) under award number R01MH105384-01A1.

Author Contributions

Conceptualization: Amir Hossein Yazdavar.

Data curation: Amir Hossein Yazdavar, John M. Meddar, Annie Myers.

Formal analysis: Amir Hossein Yazdavar.

Funding acquisition: Amir Hossein Yazdavar, Amit Sheth.

Investigation: Amir Hossein Yazdavar.

Methodology: Amir Hossein Yazdavar, Jyotishman Pathak.

Project administration: Amir Hossein Yazdavar.

Resources: Amir Hossein Yazdavar, William Romine, Pascal Hitzler.

Software: Amir Hossein Yazdavar, Mohammad Saeid Mahdavejad, Goonmeet Bajaj.

Supervision: Amir Hossein Yazdavar, Amir Hassan Monadjemi, Pascal Hitzler.

Validation: Amir Hossein Yazdavar, Pascal Hitzler.

Visualization: Amir Hossein Yazdavar.

Writing – original draft: Amir Hossein Yazdavar, Krishnaprasad Thirunarayan.

Writing – review & editing: Amir Hossein Yazdavar, Pascal Hitzler.

References

1. NIMH. How Psychotherapy and Other Treatments Can Help People Recover; 2014. Available from: <https://www.apa.org/topics/depression/recover>.
2. Marcus M, Yasamy MT, van Ommeren M, Chisholm D, Saxena S, et al. Depression: A global public health concern. WHO Department of Mental Health and Substance Abuse. 2012; 1:6–8.
3. Sullivan LE, Fiellin DA, O'Connor PG. The prevalence and impact of alcohol problems in major depression: a systematic review. *The American journal of medicine*. 2005; 118(4):330–341. <https://doi.org/10.1016/j.amjmed.2005.01.007> PMID: 15808128
4. Rudd MD, Berman AL, Joiner TE Jr, Nock MK, Silverman MM, Mandrusiak M, et al. Warning signs for suicide: Theory, research, and clinical applications. *Suicide and Life-Threatening Behavior*. 2006; 36(3):255–262. <https://doi.org/10.1521/suli.2006.36.3.255> PMID: 16805653
5. Culotta A. Lightweight methods to estimate influenza rates and alcohol sales volume from Twitter messages. *Language resources and evaluation*. 2013;.
6. Bosley JC, Zhao NW, Hill S, Shofer FS, Asch DA, Becker LB, et al. Decoding twitter: Surveillance and trends for cardiac arrest and resuscitation communication. *Resuscitation*. 2013;.
7. Yazdavar AH, Mahdavejad MS, Bajaj G, Thirunarayan K, Pathak J, Sheth A. Mental Health Analysis Via Social Media Data. In: 2018 IEEE International Conference on Healthcare Informatics (ICHI); 2018. p. 459–460.
8. Yazdavar AH, Al-Olimat HS, Ebrahimi M, Bajaj G, Banerjee T, Thirunarayan K, et al. Semi-supervised approach to monitoring clinical depressive symptoms in social media. In: Proceedings of the 2017 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining 2017. ACM; 2017. p. 1191–1198.
9. Andalibi N, Haimson OL, De Choudhury M, Forte A. Understanding social media disclosures of sexual abuse through the lenses of support seeking and anonymity. In: Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems. ACM; 2016. p. 3906–3918.
10. De Choudhury M, Gamon M, Counts S, Horvitz E. Predicting Depression via Social Media. In: ICWSM;.

11. De Choudhury M, Kiciman E, Dredze M, Coppersmith G, Kumar M. Discovering shifts to suicidal ideation from mental health content in social media. In: Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems. ACM; 2016. p. 2098–2110.
12. Manikonda L, De Choudhury M. Modeling and Understanding Visual Attributes of Mental Health Disclosures in Social Media. In: Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems. ACM; 2017. p. 170–181.
13. Andalibi N, Öztürk P, Forte A. Sensitive Self-disclosures, Responses, and Social Support on Instagram: The Case of # Depression. In: CSCW; 2017. p. 1485–1500.
14. Reece AG, Danforth CM. Instagram photos reveal predictive markers of depression. *EPJ Data Science*. 2017; 6(1):15. <https://doi.org/10.1140/epjds/s13688-017-0118-4>
15. Ahsan U, De Choudhury M, Essa I. Towards using visual attributes to infer image sentiment of social events. In: Neural Networks (IJCNN), 2017 International Joint Conference on. IEEE; 2017. p. 1372–1379.
16. MARANGA P. Social Photos Generate More Engagement: New Research; 2014. Available from: <https://www.socialmediaexaminer.com/photos-generate-engagement-research/>.
17. Cooper BB. 10 Surprising New Twitter Stats to Help You Reach More Followers; 2016. Available from: <https://blog.bufferapp.com/10-new-twitter-stats-twitter-statistics-to-help-you-reach-your-followers>.
18. Taylor M. New research: Twitter users love to watch, discover and engage with video; 2015. Available from: https://blog.twitter.com/marketing/en_us/a/2015/new-research-twitter-users-love-to-watch-discover-and-engage-with-video.html.
19. Seabrook EM, Kern ML, Rickard NS. Social networking sites, depression, and anxiety: a systematic review. *JMIR mental health*. 2016; 3(4). <https://doi.org/10.2196/mental.5842>
20. Liu L, Preotiu-Pietro D, Samani ZR, Moghaddam ME, Ungar LH. Analyzing Personality through Social Media Profile Picture Choice. In: ICWSM; 2016. p. 211–220.
21. Montesano A, Feixas G, Caspar F, Winter D. Depression and Identity: Are Self-Constructions Negative or Conflicting? *Frontiers in psychology*. 2017; 8:877. <https://doi.org/10.3389/fpsyg.2017.00877> PMID: 28611716
22. Nolen-Hoeksema S. Sex differences in unipolar depression: evidence and theory. *Psychological bulletin*. 1987; 101(2):259. <https://doi.org/10.1037/0033-2909.101.2.259> PMID: 3562707
23. McManus S, Bebbington P, Jenkins R, Brugha T. Mental Health and Wellbeing in England: Adult Psychiatric Morbidity Survey 2014: a Survey Carried Out for NHS Digital by NatCen Social Research and the Department of Health Sciences, University of Leicester. NHS Digital; 2016.
24. Angst J, Gamma A, Gastpar M, Lépine JP, Mendlewicz J, Tylee A. Gender differences in depression. *European archives of psychiatry and clinical neuroscience*. 2002; 252(5):201–209. <https://doi.org/10.1007/s00406-002-0381-6> PMID: 12451460
25. Meltzer H, Gill B, Petticrew M. The prevalence of psychiatric morbidity among adults living in private households. In: The prevalence of psychiatric morbidity among adults living in private households; 1995.
26. Cook MN, Peterson J, Sheldon C. Adolescent depression: an update and guide to clinical decision making. *Psychiatry (Edgmont)*. 2009; 6(9):17.
27. Nolen-Hoeksema S, Girgus JS. The emergence of gender differences in depression during adolescence. *Psychological bulletin*. 1994; 115(3):424. <https://doi.org/10.1037/0033-2909.115.3.424> PMID: 8016286
28. Ruch DA, Sheftall AH, Schlagbaum P, Rausch J, Campo JV, Bridge JA. Trends in suicide among youth aged 10 to 19 years in the United States, 1975 to 2016. *JAMA network open*. 2019; 2(5): e193886–e193886. <https://doi.org/10.1001/jamanetworkopen.2019.3886> PMID: 31099867
29. Levi G, Hassner T. Age and gender classification using convolutional neural networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops; 2015. p. 34–42.
30. De Choudhury M, Counts S, Horvitz E. Social media as a measurement tool of depression in populations. In: Proceedings of the 5th Annual ACM Web Science Conference. ACM; 2013. p. 47–56.
31. Nguyen T, Phung D, Dao B, Venkatesh S, Berk M. Affective and content analysis of online depression communities. *IEEE Transactions on Affective Computing*. 2014; 5(3):217–226. <https://doi.org/10.1109/TAFFC.2014.2315623>
32. Coppersmith G, Dredze M, Harman C, Hollingshead K, Mitchell M. CLPsych 2015 shared task: Depression and PTSD on Twitter. In: Proceedings of the 2nd Workshop on Computational Linguistics and Clinical Psychology; 2015.
33. Resnik P, Armstrong W, Claudino L, Nguyen T, Nguyen VA, Boyd-Graber J. Beyond LDA: exploring supervised topic modeling for depression-related language in Twitter. In: Proceedings of the 2nd

- Workshop on Computational Linguistics and Clinical Psychology: From Linguistic Signal to Clinical Reality; 2015.
34. Preotiuc-Pietro D, Eichstaedt J, Park G, Sap M, Smith L, Tobolsky V, et al. The role of personality, age and gender in tweeting about mental illnesses. In: NAACL HLT; 2015.
 35. Yates A, Cohan A, Goharian N. Depression and Self-Harm Risk Assessment in Online Forums. arXiv preprint arXiv:170901848. 2017;.
 36. Milne DN, Pink G, Hachey B, Calvo RA. CLPsych 2016 Shared Task: Triaging content in online peer-support forums. In: Proceedings of the Third Workshop on Computational Linguistics; 2016.
 37. Wang Y, Wang S, Tang J, Liu H, Li B. Unsupervised Sentiment Analysis for Social Media Images. In: IJCAI; 2015. p. 2378–2379.
 38. Abdullah S, Murnane EL, Costa JM, Choudhury T. Collective smile: Measuring societal happiness from geolocated images. In: Proceedings of the 18th ACM Conference on Computer Supported Cooperative Work & Social Computing. ACM; 2015. p. 361–374.
 39. Garimella VRK, Alfayad A, Weber I. Social media image analysis for public health. In: Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems. ACM; 2016. p. 5543–5547.
 40. Duong CT, Lebre R, Aberer K. Multimodal Classification for Analysing Social Media. arXiv preprint arXiv:170802099. 2017;.
 41. Shen G, Jia J, Nie L, Feng F, Zhang C, Hu T, et al. Depression detection via harvesting social media: A multimodal dictionary learning solution. In: Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence (IJCAI-17); 2017. p. 3838–3844.
 42. Mislove A, Lehmann S, Ahn YY, Onnela JP, Rosenquist JN. Understanding the Demographics of Twitter Users. ICWSM. 2011; 11:5th.
 43. Ebrahimi M, Ebrahimi M, Yazdavar AH, Yazdavar AH, Salim N, Salim N, et al. Recognition of side effects as implicit-opinion words in drug reviews. *Online Information Review*. 2016; 40(7):1018–1032. <https://doi.org/10.1108/OIR-06-2015-0208>
 44. Yazdavar AH, Ebrahimi M, Salim N. Fuzzy based implicit sentiment analysis on quantitative sentences. arXiv preprint arXiv:170100798. 2017;.
 45. Wakamiya S, Matsune S, Okubo K, Aramaki E. Causal Relationships Among Pollen Counts, Tweet Numbers, and Patient Numbers for Seasonal Allergic Rhinitis Surveillance: Retrospective Analysis. *Journal of medical Internet research*. 2019; 21(2):e10450. <https://doi.org/10.2196/10450> PMID: 30785411
 46. Zagheni E, Garimella VRK, Weber I, et al. Inferring international and internal migration patterns from twitter data. In: Proceedings of the 23rd International Conference on World Wide Web. ACM; 2014. p. 439–444.
 47. Alowibdi JS, Buy UA, Yu P. Language independent gender classification on Twitter. In: Proceedings of the 2013 IEEE/ACM international conference on advances in social networks analysis and mining. ACM; 2013. p. 739–743.
 48. Mueller J, Stumme G. Gender inference using statistical name characteristics in twitter. In: Proceedings of the The 3rd Multidisciplinary International Social Networks Conference on SocialInformatics 2016, Data Science 2016. ACM; 2016. p. 47.
 49. An J, Weber I. # greysanatomy vs. # yankees: Demographics and Hashtag Use on Twitter. In: Tenth International AAAI Conference on Web and Social Media; 2016.
 50. Kosinski M, Stillwell D, Graepel T. Private traits and attributes are predictable from digital records of human behavior. *Proceedings of the National Academy of Sciences*. 2013; 110(15):5802–5805. <https://doi.org/10.1073/pnas.1218772110>
 51. Burger JD, Henderson J, Kim G, Zarrella G. Discriminating gender on Twitter. In: Proceedings of the conference on empirical methods in natural language processing. Association for Computational Linguistics; 2011. p. 1301–1309.
 52. Rosenthal S, McKeown K. Age prediction in blogs: A study of style, content, and online behavior in pre-and post-social media generations. In: Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies-Volume 1. Association for Computational Linguistics; 2011. p. 763–772.
 53. Nguyen D, Gravel R, Trieschnigg D, Meder T. "How Old Do You Think I Am?" A Study of Language and Age in Twitter. In: ICWSM; 2013.
 54. Sloan L, Morgan J, Burnap P, Williams M. Who tweets? Deriving the demographic characteristics of age, occupation and social class from Twitter user meta-data. *PloS one*. 2015; 10(3):e0115545. <https://doi.org/10.1371/journal.pone.0115545> PMID: 25729900
 55. Standard Occupational Classification;. Available from: <https://www.bls.gov/soc/>.

56. Zhang J, Hu X, Zhang Y, Liu H. Your Age Is No Secret: Inferring Microbloggers' Ages via Content and Interaction Analysis. In: ICWSM; 2016. p. 476–485.
57. Nguyen D, Trieschnigg D, Dođruöz AS, Gravel R, Theune M, Meder T, et al. Why gender and age prediction from tweets is hard: Lessons from a crowdsourcing experiment. In: Proceedings of COLING 2014, the 25th International Conference on Computational Linguistics: Technical Papers; 2014. p. 1950–1961.
58. Masi I, Tran AT, Hassner T, Leksut JT, Medioni G. Do we really need to collect millions of faces for effective face recognition? In: European Conference on Computer Vision. Springer; 2016. p. 579–596.
59. Barrick CB, Taylor D, Correa EI. Color sensitivity and mood disorders: biology or metaphor? *Journal of affective disorders*. 2002; 68(1):67–71. [https://doi.org/10.1016/s0165-0327\(00\)00358-x](https://doi.org/10.1016/s0165-0327(00)00358-x) PMID: 11869784
60. Lüscher M. The Luscher color test. Simon and Schuster; 1990.
61. Nikolaenko N. Role of the dominant and nondominant hemispheres in the perception and naming of color. *Human physiology*. 1981;.
62. Fernandes TMP, Andrade SM, de Andrade MJO, Nogueira RMTBL, Santos NA. Colour discrimination thresholds in type 1 Bipolar Disorder: a pilot study. *Scientific reports*. 2017; 7(1):16405. <https://doi.org/10.1038/s41598-017-16752-0> PMID: 29180712
63. Carruthers HR, Morris J, Tarrier N, Whorwell PJ. The Manchester Color Wheel: development of a novel way of identifying color choice and its validation in healthy, anxious and depressed individuals. *BMC medical research methodology*. 2010; 10(1):12. <https://doi.org/10.1186/1471-2288-10-12> PMID: 20144203
64. Nie J, Wei Z, Li Z, Yan Y, Huang L. Understanding personality of portrait by social embedding visual features. *Multimedia Tools and Applications*. 2019; 78(1):727–746. <https://doi.org/10.1007/s11042-017-5577-x>
65. Mejova Y, Haddadi H, Noulas A, Weber I. # foodporn: Obesity patterns in culinary interactions. In: Proceedings of the 5th international conference on digital health 2015. ACM; 2015. p. 51–58.
66. Hassanpour S, Tomita N, DeLise T, Crosier B, Marsch LA. Identifying substance use risk based on deep neural networks and Instagram social media data. *Neuropsychopharmacology*. 2019; 44(3):487. <https://doi.org/10.1038/s41386-018-0247-x> PMID: 30356094
67. Torous J, Larsen ME, Depp C, Cosco TD, Barnett I, Nock MK, et al. Smartphones, sensors, and machine learning to advance real-time prediction and interventions for suicide prevention: a review of current progress and next steps. *Current psychiatry reports*. 2018; 20(7):51. <https://doi.org/10.1007/s11920-018-0914-y> PMID: 29956120
68. Celli F, Bruni E, Lepri B. Automatic personality and interaction style recognition from facebook profile pictures. In: Proceedings of the 22nd ACM international conference on Multimedia. ACM; 2014. p. 1101–1104.
69. Segalin C, Celli F, Polonio L, Kosinski M, Stillwell D, Sebe N, et al. What your Facebook profile picture reveals about your personality. In: Proceedings of the 25th ACM international conference on Multimedia. ACM; 2017. p. 460–468.
70. Jourard SM. Self-disclosure: An experimental analysis of the transparent self. 1971;.
71. Twitter API; Available from: <https://developer.twitter.com/en/docs/tweets/search/api-reference/get-search-tweets.html#>.
72. Depression Detector; Available from: <https://github.com/halolimat/Depression-Detector>.
73. Liao L, Jiang J, Lim EP, Huang H. A study of age gaps between online friends. In: Proceedings of the 25th ACM conference on Hypertext and social media. ACM; 2014. p. 98–106.
74. Duggan M, Ellison NB, Lampe C, Lenhart A, Madden M. Demographics of key social networking platforms. Pew Research Center. 2015; 9.
75. Al Zamil F, Liu W, Ruths D. Homophily and Latent Attribute Inference: Inferring Latent Attributes of Twitter Users from Neighbors. ICWSM. 2012; 270:2012.
76. Ford ES, Giles WH, Dietz WH. Prevalence of the metabolic syndrome among US adults: findings from the third National Health and Nutrition Examination Survey. *Jama*. 2002; 287(3):356–359. <https://doi.org/10.1001/jama.287.3.356> PMID: 11790215
77. Barney LJ, Griffiths KM, Jorm AF, Christensen H. Stigma about depression and its impact on help-seeking intentions. *Australian & New Zealand Journal of Psychiatry*. 2006; 40(1):51–54. <https://doi.org/10.1080/j.1440-1614.2006.01741.x>
78. Zhou E, Fan H, Cao Z, Jiang Y, Yin Q. Extensive facial landmark localization with coarse-to-fine convolutional network cascade. In: Proceedings of the IEEE International Conference on Computer Vision Workshops; 2013. p. 386–391.

79. Face ++. Available from: <https://www.faceplusplus.com>.
80. Emotion classification;. Available from: https://en.wikipedia.org/wiki/Emotion_classification.
81. Datta R, Joshi D, Li J, Wang JZ. Studying aesthetics in photographic images using a computational approach. In: European Conference on Computer Vision. Springer; 2006. p. 288–301.
82. Huang KQ, Wang Q, Wu ZY. Natural color image enhancement and evaluation algorithm based on human visual system. *Computer Vision and Image Understanding*. 2006; 103(1):52–63.
83. San Pedro J, Siersdorfer S. Ranking and classifying attractiveness of photos in folksonomies. In: Proceedings of the 18th international conference on World wide web. ACM; 2009. p. 771–780.
84. Python-tesseract: an optical character recognition (OCR) tool for python;. Available from: <https://pypi.org/project/pytesseract/>.
85. Ebrahimi M, Yazdavar AH, Sheth A. Challenges of sentiment analysis for dynamic events. *IEEE Intelligent Systems*. 2017; 32(5):70–75. <https://doi.org/10.1109/MIS.2017.3711649>
86. How the words we use in everyday language reveal our thoughts, feelings, personality, and motivations;. Available from: <http://iwc.wpengine.com/>.
87. Schwartz HA, Eichstaedt JC, Kern ML, Dziurzynski L, Ramones SM, Agrawal M, et al. Personality, gender, and age in the language of social media: The open-vocabulary approach. *PLoS one*. 2013; 8(9):e73791. <https://doi.org/10.1371/journal.pone.0073791> PMID: 24086296
88. Chung C, Pennebaker JW. The psychological functions of function words. *Social communication*. 2007; 1:343–359.
89. Mukherjee A, Liu B. Improving gender classification of blog authors. In: Proceedings of the 2010 conference on Empirical Methods in natural Language Processing. Association for Computational Linguistics; 2010. p. 207–217.
90. Argamon S, Koppel M, Pennebaker JW, Schler J. Mining the blogosphere: Age, gender and the varieties of self-expression. *First Monday*. 2007; 12(9). <https://doi.org/10.5210/firstmonday.v12i9.2003>
91. Newman ML, Groom CJ, Handelman LD, Pennebaker JW. Gender differences in language use: An analysis of 14,000 text samples. *Discourse Processes*. 2008; 45(3):211–236. <https://doi.org/10.1080/01638530802073712>
92. Pennebaker JW, Stone LD. Words of wisdom: Language use over the life span. *Journal of personality and social psychology*. 2003; 85(2):291. <https://doi.org/10.1037/0022-3514.85.2.291> PMID: 12916571
93. Rude S, Gortner EM, Pennebaker J. Language use of depressed and depression-vulnerable college students. *Cognition & Emotion*. 2004; 18(8):1121–1133. <https://doi.org/10.1080/02699930441000030>
94. Theories of Adult Development;. Available from: <https://study.com/academy/lesson/theories-of-adult-development-levinson-valliant-neugarten.html>.
95. Kintgen-Andrews J. Critical thinking and nursing education: Perplexities and insights. *Journal of Nursing Education*. 1991; 30(4):152–157. PMID: 1646306
96. Critical Thinking and the Three Stages of Cognitive Development;. Available from: <https://creativityandcriticalthinking.wordpress.com/the-evolution-from-pre-k-to-college/critical-thinking-and-the-three-stages-of-cognitive-development/>.
97. Eckert P. Age as a sociolinguistic variable. *The handbook of sociolinguistics*. 2017; p. 151–167.
98. Newman ML, Pennebaker JW, Berry DS, Richards JM. Lying words: Predicting deception from linguistic styles. *Personality and social psychology bulletin*. 2003; 29(5):665–675. <https://doi.org/10.1177/0146167203029005010> PMID: 15272998
99. Lies Can Point to Mental Disorders or Signal Normal Growth;. Available from: <https://www.nytimes.com/1988/05/17/science/lies-can-point-to-mental-disorders-or-signal-normal-growth.html>.
100. Orth U, Erol RY, Luciano EC. Development of self-esteem from age 4 to 94 years: A meta-analysis of longitudinal studies. *Psychological bulletin*. 2018; 144(10):1045. <https://doi.org/10.1037/bul0000161> PMID: 30010349
101. PHQ-9;. Available from: https://www.phqscreeners.com/sites/g/files/q10049256/f/201412/PHQ-9_English.pdf.
102. Kaye BK, Sapolsky BS. Watch your mouth! An analysis of profanity uttered by children on prime-time television. *Mass Communication & Society*. 2004; 7(4):429–452. https://doi.org/10.1207/s15327825mcs0704_4
103. The Surprising Health Benefits of Swearing;. Available from: <https://psychcentral.com/blog/the-surprising-health-benefits-of-swearing/>.
104. Aging and Male Sexual Desire II: Physical Factors;. Available from: <https://www.psychologytoday.com/us/blog/mindful-sex/201301/aging-and-male-sexual-desire-ii-physical-factors>.

105. Sap M, Park G, Eichstaedt J, Kern M, Stillwell D, Kosinski M, et al. Developing age and gender predictive lexica over social media. In: Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP); 2014. p. 1146–1151.
106. Snoek CG, Worring M, Smeulders AW. Early versus late fusion in semantic video analysis. In: Proceedings of the 13th annual ACM international conference on Multimedia. ACM; 2005. p. 399–402.
107. Kursa MB, Rudnicki WR, et al. Feature selection with the Boruta package. *J Stat Softw.* 2010; 36(11):1–13. <https://doi.org/10.18637/jss.v036.i11>
108. Chen T, Guestrin C. Xgboost: A scalable tree boosting system. In: Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining. ACM; 2016. p. 785–794.
109. XGBOOSTExplainer;. Available from: <https://github.com/AppliedDataSciencePartners/xgboostExplainer>.
110. Nadeem M. Identifying depression on Twitter. arXiv preprint arXiv:160707384. 2016;.
111. Coppersmith G, Ngo K, Leary R, Wood A. Exploratory analysis of social media prior to a suicide attempt. In: Proceedings of the Third Workshop on Computational Linguistics and Clinical Psychology; 2016.
112. Coppersmith G, Dredze M, Harman C. Quantifying mental health signals in Twitter. In: Proceedings of the Workshop on Computational Linguistics and Clinical Psychology: From Linguistic Signal to Clinical Reality; 2014. p. 51–60.
113. Preotjuc-Pietro D, Eichstaedt J, Park G, Sap M, Smith L, Tobolsky V, et al. The role of personality, age, and gender in tweeting about mental illness. In: Proceedings of the 2nd workshop on computational linguistics and clinical psychology: From linguistic signal to clinical reality; 2015. p. 21–30.
114. Tsugawa S, Kikuchi Y, Kishino F, Nakajima K, Itoh Y, Ohsaki H. Recognizing depression from twitter activity. In: Proceedings of the 33rd Annual ACM Conference on Human Factors in Computing Systems; 2015.
115. Wang P, Xu J, Xu B, Liu CL, Zhang H, Wang F, et al. Semantic Clustering and Convolutional Neural Network for Short Text Categorization. In: *ACL (2)*; 2015. p. 352–357.
116. Word2vec;. Available from: <https://github.com/loretoparisi/word2vec-twitter>.
117. Mikolov T, Sutskever I, Chen K, Corrado GS, Dean J. Distributed representations of words and phrases and their compositionality. In: *Advances in neural information processing systems*; 2013. p. 3111–3119.
118. Gunaratna K, Yazdavar AH, Thirunarayan K, Sheth A, Cheng G. Relatedness-based multi-entity summarization. In: *IJCAI: proceedings of the conference.* vol. 2017. NIH Public Access; 2017. p. 1060.

Semi-Supervised Approach to Monitoring Clinical Depressive Symptoms in Social Media

Amir Hossein Yazdavar*, Hussein S. Al-Olimat*, Monireh Ebrahimi*, Goonmeet Bajaj*, Tanvi Banerjee*,
Krishnaprasad Thirunarayan*, Jyotishman Pathak** and Amit Sheth *

* Kno.e.sis Center, Wright State University, Dayton, OH, USA

Email: {amir;hussein;monireh;goonmeet;tanvi;tkprasad;amit}@knoesis.org

** Division of Health Informatics, Cornell University, New York, NY, USA

Email: jyp2001@med.cornell.edu

Abstract—With the rise of social media, millions of people are routinely expressing their moods, feelings, and daily struggles with mental health issues on social media platforms like Twitter. Unlike traditional observational cohort studies conducted through questionnaires and self-reported surveys, we explore the reliable detection of clinical depression from tweets obtained unobtrusively. Based on the analysis of tweets crawled from users with self-reported depressive symptoms in their Twitter profiles, we demonstrate the potential for detecting clinical depression symptoms which emulate the PHQ-9 questionnaire clinicians use today. Our study uses a semi-supervised statistical model to evaluate how the duration of these symptoms and their expression on Twitter (in terms of word usage patterns and topical preferences) align with the medical findings reported via the PHQ-9. Our proactive and automatic screening tool is able to identify clinical depressive symptoms with an accuracy of 68% and precision of 72%.

Keywords—Semi-supervised Machine Learning, Natural Language Processing, Social Media, Mental Health

I. INTRODUCTION

A common global effort to manage depression involves detecting depression through survey-based methods via phone or online questionnaires¹. However, these studies suffer from underrepresentation, sampling biases, and incomplete information. Additionally, large temporal gaps between data collection and the dissemination of findings can delay the administration of timely and appropriate remedial measures. Cognitive bias, which prevents participants from giving truthful responses, is yet another limitation [1]. In contrast, Twitter is a valuable resource for learning about users' feelings, emotions, behaviors, and decisions that reflect their mental health as they are experiencing the ups and the downs in real-time. For example, news headlines such as "Twitter fail: Teen Sent 144 Tweets Before Committing Suicide & No One Helped" and "Jim Carrey's Girlfriend: Her Last Tweet Before Committing Suicide 'Signing Off'", illustrate the expression of emotional

¹<https://www.cdc.gov/mentalhealthsurveillance/>

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

ASONAM '17, July 31 - August 03, 2017, Sydney, Australia

© 2017 Association for Computing Machinery.

ACM ISBN 978-1-4503-4993-2/17/07?/\$15.00

<http://dx.doi.org/10.1145/3110025.3123028>

turmoils in tweets and subsequent deliberate actions in the physical world.

In recent years, much progress has been made in studying mood and mental health through social media content [2; 3; 4; 5; 6]. These studies can be categorized into two major groups; namely, lexicon-based [7; 8], and supervised [5; 9; 10; 11; 12]. These studies suggested the individual's language style, emotion, ego-network, and user engagement as discriminating features to recognize depression-indicative posts. However, the lexicon-based approaches suffer from low recall and are highly dependent on the quality of the created lexicon. On the other hand, supervised approaches require labor intensive annotation of a huge dataset. Besides, suffering from clinical depression is more than feeling down for a few days [13]. Indeed, clinical depression is diagnosed through a set of predefined symptoms which last for a fixed period.

Inspired by that, we develop a statistical model which emulates traditional observational cohort studies conducted through online questionnaires by extracting, categorizing, and unobtrusively monitoring different symptoms of depression by modeling user-generated content in social media as a mixture of underlying topics evolving over time. To our knowledge, this is the first study that incorporates temporal analysis of user-generated content on social media for capturing these tell-tale symptoms. We crawled 23 million tweets posted by over 45,000 Twitter users who self-reported symptoms of depression in their profile descriptions.

The present study answers the following questions: 1) How well can textual content in social media be harnessed to reliably capture a user's symptoms of clinical depression over time and build a proactive and automatic depression screening tool? 2) Are there any underlying common themes among depressed users?

We assess the level of depression expressed in tweets for each user profile in our dataset by integrating a lexicon-based method (top-down processing) with a data-driven method (bottom-up processing). Leveraging the clinical articulation of depression, we build a depression lexicon that contains common depression symptoms from the established clinical assessment questionnaire PHQ-9 [13]. We rank the terms and compile a list of informative lexicon terms for each user and use them as seed terms to discover latent topics (depression symptoms) discussed by the subject in his/her tweets (bottom-up processing). We develop a probabilistic topic modeling

over user tweets with partial supervision (by leveraging seeded clusters), named *semi-supervised topic modeling over time* (ssToT), to monitor clinical depression symptoms. We apply ssToT to derive the per user topic (depression symptoms) distribution and per topic word distribution to screen and determine a trend of symptoms over time.

The major contributions of this multidisciplinary study, conducted by a team of computer scientists and mental health experts are two-fold: **First**, we create a lexicon of depression symptoms which are likely to appear in the generated content of depressed individuals; **second**, we develop a semi-supervised statistical model to extract, categorize, and monitor depression symptoms for continuous temporal analysis of an individual's tweets. Empirical evaluations show our model is superior to five baselines in terms of the quality of learned topics (clinical depression symptoms).

II. RELATED WORK

Several efforts have attempted to automatically detect depression in social media content using machine learning approach. Conducting a retrospective study over the content of tweets posted by depressed and non-depressed individuals for one year, [14] characterize depression based on factors such as language, emotion, style, ego-network, and user engagement. They utilize these distinguishing characteristics to build a classifier to predict the likelihood of depression in a post [14] or in an individual [9]. In another study, [10] leverage affective aspect, linguistic style, and topics as a feature for detecting depressed communities. [15] employ various features including sentence polarity for detection of depressed individuals in Twitter. They use n-grams, psychological categories, and emoticons as features. Similarly, [16] apply a classifier for identifying potential cases of social network mental disorders. Moreover, there have been significant advances in the field by introducing a shared task [11] at the Computational Linguistics and Clinical Psychology Workshop (CLP 2015) focusing on methods for identifying depressed users on Twitter. A corpus of nearly 1,800 Twitter users was built for evaluation; among the several participants in the shared task, the best models employed topic modeling [2], and various features such as bag of words, LIWC features, metadata and clustering features [17]. Another related line of research focused on capturing suicide and self-harm signals from Twitter posts [18]. Through analysis of tweets posted by individuals attempting committing suicide, [19] indicate quantifiable signals of suicidal ideations. Moreover, the 2016 ACL Computational Linguistics and Clinical Psychology Workshop [20] defined a shared task on detecting the severity of the mental health forum posts. *All* of these studies define some discriminative features to classify depression in user-generated content in a message, for a user or at a community level. However, our approach facilitates detection of depression through fine-grained temporal monitoring of subjects' behavior by analyzing the symptoms of depression mirrored in their topics of interest and word usage. Apart from that, what makes our model different is that it does not require any labeled dataset.

In the context of *lexicon-based* approaches, [8] use a dictionary-based method for assigning an overall depression score to subjects. They count all the phrases that are matched

with depression indicators without considering separate symptom categories. [21] study the usage of keyword the "depression" in tweets. They find initial evidence that individuals tweet about their depression and even disclose updates about their mental health treatment on Twitter. They found an association between excessive use of negative-emotions-related words and having a major depressive disorder. In contrast, no relation has been found in the use of positive-emotions-related words and depression. Similarly, [7] propose an NLP methodology for automatic screening of depression using a depression lexicon incorporating both metaphorical and non-metaphorical words and phrases. They perform a web search to retrieve documents containing "depression is like *" pattern. However, in natural language, words can be ambiguous. For instance, depression may be used to express different concepts such as "economic depression", "great depression", "depression era", and "tropical depression". Moreover, neurotypical people use this term to express their transient sadness. For instance, consider: "I am depressed, I have a final exam tomorrow". Furthermore, the experience of depression may be expressed implicitly, making a lexicon-based approach insufficient for accurate fine-grained analysis of depression symptoms over time. Another inherent drawback of all lexicon-based methods is their high precision at the expense of low recall and lack of context-sensitivity. For instance, "...sleep forever..." may indicate suicidal thoughts rather than the act of sleeping. In short, our study differs from existing works in that we developed a statistical model for the linguistic analysis of social media content authored by a subject by seeking depression indicators and their variation over time.

III. PROPOSED APPROACH

The Diagnostic and Statistical Manual of Mental Disorders (DSM)² suggests that clinical depression can be diagnosed through the presence of a set of symptoms over a fixed period of time. The PHQ-9³ is a nine item depression scale, which incorporates DSM-V. It can be utilized to screen, diagnose, and measure the severity of depression. Our research hypothesis is that depressed individuals discuss their symptoms on Twitter. Symptoms of depression include decreased pleasure in most activities (S1), feeling down (S2), sleep disorders (S3), loss of energy (S4), a significant change in appetite (S5), feeling worthless (S6), concentration problems (S7), hyper/lower activity (S8), and suicidal thoughts (S9). This is a top-down definition of depressive disorder through its "symptomatology". To validate this hypothesis, we first manually examined symptoms in a random selection of 100 user profiles in our dataset. Table I illustrates a sample of anonymized tweets and their associated symptoms in PHQ-9. This table highlights the importance of developing appropriate models of textual content to capture depressive behavior on Twitter.

Motivated by these observations, we investigate two approaches for detecting symptoms of clinical depression on Twitter, emulating the PHQ-9 questionnaire. The first approach captures clinical depression using bottom-up processing of user tweets and distributional semantics to uncover symptoms of depression via related word clusters. The second approach hybridizes the first approach with top-down processing by

²<https://www.psychiatry.org/psychiatrists/practice/dsm/dsm5>

³http://www.cqaimh.org/pdf/tool_phq9.pdf

TABLE I. TWEETS SAMPLE AND THEIR ASSOCIATED SYMPTOMS

PHQ-9 Symptoms	Short-text Document
Lack of Interest	I've not replied all day due to total lack of interest, depressed probs
Feeling Down	i feel like i'm falling apart.
Sleep Disorder	Night guys. Hope you sleep better than me.
Lack of Energy	so tired, so drained, so done
Eating Disorder	I just wanna be skinny and beautiful
Low Self-esteem	I am disgusted with myself.
Concentration Problems	I couldn't concentrate to classes at all can't stop thinking
Hyper/Lower Activity	so stressed out I cant do anything
Suicidal Thoughts	I want summer but then i don't... It'll be harder to hide my cuts.

using the lexicon terms to guide the extraction of symptoms from tweets.

A. Bottom-up processing: LDA

In health data mining, the problem of discovering latent topics represents a promising research area [22]. We hypothesize that by analyzing a user's topic preferences (what) and word usage (how) we can monitor depression symptoms. Our approach is based on latent variable topic models, more specifically, Latent Dirichlet Allocation (LDA). LDA is an unsupervised method that views a document as a mixture of latent topics, where a topic is a distribution of co-occurring words. Different terms expressing a related facet would be grouped together under the same topic. We apply LDA to extract latent topics discussed by users in our dataset.

Not surprisingly, the topics learned by LDA are not granular and specific enough to correspond to depressive symptoms. Several prior studies also highlight that the results from the traditional LDA do not correlate well with human judgments [23; 24]. Some work has been done to guide the discovery of latent topics in LDA by incorporating domain knowledge in different ways; from defining a set of First-Order Logic (FOL) rules [25] to constraining the occurrence of some terms together by encoding a set of Must-Links and Cannot-Links associated with the domain knowledge [26] or in the context of aspect extraction for sentiment analysis by providing some relevant terms for a few aspects [27].

The key difference between our seeding model and this study is that we supervise the topics at the token level rather than measure the distribution over a predefined list of terms. In particular, we restrict the occurrence of relevant tokens within the specified topics. We will further explain the seeding approach in the following section.

B. Hybrid processing: Proposed ssToT Model

The basis of traditional LDA is the frequency of the co-occurrence of terms in various contexts. This syntactic approach often results in many terms from different symptom categories being merged into a single topic. By constraining symptom-related seed terms so that they only appear in a single topic, we bias the "bottom-up" learned topics to align with expected "top-down" symptom categories. In particular, we add supervision to LDA, by using terms that are strongly related to the 9 depression symptoms as seeds of the topical clusters and guide the model to aggregate semantically-related terms into the same cluster.

To generate a set of seed terms for each symptom category, we leverage the lexicon as background knowledge. In particular, in collaboration with our psychologist clinician, we built a lexicon of depression-related terms that are likely to be utilized by individuals suffering from depression. We use patient health questionnaire (PHQ-9) categories as a predefined list of depression symptoms. Furthermore, given the colloquial language of social media, we use Urban Dictionary (a crowd-sourced online dictionary of slang words and phrases)⁴ for expanding the lexicon using the synset of each of the nine PHQ-9 depression symptoms categories. We also employ Big Huge Thesaurus⁵ to obtain synonyms for each symptom category. The consistency of the built lexicon with psychologist's requirements has been vetted by our psychologist collaborators. After several rounds of refinement by domain experts, the final lexicon contains more than 1,620 depression-related symptoms categorized into nine different clinical depression symptom categories which are likely to appear in the tweets of individuals suffering from clinical depression.

However, there are important challenges to overcome in order to effectively leverage our lexicon for compiling a seed cluster. First, social media users often use diverse terms to express a specific concept. They use creative descriptive metaphorical phrases and explanations for symptoms. One may say, "I'm so exhausted all time" while another may say "so tired, so drained, so done" while both of these utterances discuss the unique medical concept "Lack of Energy". Second, language of social media contains polysemous words in its vocabulary. Their interpretation requires context for Word Sense Disambiguation (WSD). For instance, "Cut my finger opening a can of fruit" and "scars don't heal when you keep cutting" use "cut" in different contexts and senses.

To address the *first* challenge, our algorithm automatically generates a personalized set of seed terms per user which is a subset of the available terms in the lexicon. In this manner, a list of highly informative seeds will be generated per user. For the above examples, the term "exhausted" would be a seed for the first user while "drained" and "tired" would be the seeds for the second user. To address the *second* challenge, given the recent advances in sentiment analysis techniques [28; 29; 30], we disambiguate a polysemous word based on the sentiment polarity of its enclosing sentence. We include a term as a seed only if the enclosing context has negative sentiment. We perform sentiment analysis using the Python TextBlob⁶, a standard library, which determines positive/neutral/negative polarity for any document. For the above example, "cut" is not a seed for the first user, but is a seed for the second user, as the first tweet reflects a neutral sentiment while the second tweet indicates a negative sentiment.

On the other hand, experiencing clinical depression is more than feeling down for a few days. According to PHQ-9 clinical depression symptoms should persist for a few weeks. Hence, temporal monitoring of symptoms is crucial.

⁴<http://www.urbandictionary.com/>

⁵<https://words.bighugelabs.com/>

⁶<https://textblob.readthedocs.io/en/dev/>

IV. ALGORITHM

Motivated by the above observations, we propose our framework to automatically analyze user behavior by continuously monitoring their social media content over time intervals. To this end, the proposed approach enriches the LDA model’s expressiveness by introducing a predefined set of seed terms. We divide each user’s collection of preprocessed tweets into a set of tweet buckets using a specific time interval of d days. The generative process of the proposed model for a corpus C of individual user’s tweets consisting of B buckets is shown in Algorithm-1.

Algorithm 1 The generative process of ssToT

```

1: procedure ANALYZETWITTERPROFILE
2:   for each symptom (topic)  $s \in 1, 2, \dots, 9$  do
3:     Draw a distribution over terms and seed sets
4:      $\Phi_s \sim \text{Dirichlet}(\beta)$ 
5:   end for
6:   for each bucket (document)  $b \in 1, 2, \dots, B$  do
7:     Draw a distribution over topics  $\theta_b \sim \text{Dirichlet}(\alpha)$ 
8:     for each word  $w_i \in b$  do
9:       Choose a symptom (topic)  $s_i \sim \text{multinomial}(\theta_b)$ 
10:      Choose a word  $w_i \sim \text{multinomial}(\Phi_{s_i})$ 
11:     end for
12: end procedure

```

In Algorithm-1, θ shows the distribution of symptoms over buckets while Φ is the distribution of words per symptom. We employ Gibbs Sampling to approximate the posterior distribution over the assignment of words to topics, $P(s|w)$. We then estimate Φ and θ using this posterior distribution. Our strategy for discovering symptoms (topics) differs from previous methods as we incorporate prior knowledge into the inference by assigning the pre-defined seed terms into only one of the symptoms (topics). Inspired by [23], we adapt the Gibbs Sampling equation by restricting a topic s_i to a single corresponding value for each user-specific seed term or phrase. Each term w_i is assigned to the largest probability symptom associated with it in Φ . We change the probability of a symptom over a bucket to zero if the number of seed terms associated with it is less than a threshold τ . Similarly, to filter out polysemous seed terms, we aggregate the sentiment polarity of all sentences containing all seed terms over a bucket. If the aggregated polarity is positive, we assign the probability of zero to all symptoms in that bucket. Finally, we visualize the probability of each symptom over the bucket in matrix θ for further analysis and monitoring. Apart from that, if the probability of a symptom is more than a threshold τ , the symptom would be assigned to the bucket as a label. In this manner, our model can be utilized as a multi-label classifier over a time interval. The quality of our multi-label classifier is evaluated as follows.

V. EXPERIMENTAL RESULTS

We first discuss data collection procedure, followed by qualitative and quantitative analysis. We highlight that since this study analyzes individual’s behavioral health information, which may be considered as sensitive, in our datasets, we anonymized users’ real identities as per the approved institutional review board (IRB) protocol.

Dataset: We created a dataset containing 45,000 Twitter users who self-declared their depression and 2,000 “undeclared” users who were collected randomly. In particular,

for collecting self-declared depressed individual’s profiles, we utilize a subset of highly informative depressive indicative terms in our lexicon and find the profiles that contain these terms in their description. Afterwards, we crawled the tweets, the tweets’ timestamp, and the list of friends and followers of these users. After removing the profiles with less than 100 tweets, we obtained 7,046 users with 21 million timestamped tweets, with each user contributing at most 3,200 tweets due to the Twitter Search API limitation. Next, we randomly sampled 2,000 profiles of users with self-reported depression symptoms and 2,000 random users who do not have any depression terms in their profile descriptions. We denote this subset of 4,000 users by U . We preprocess these tweets by changing the space delimiter into underscore in all phrases in the tweets that are listed in our lexicon as a seed phrase (e.g., lack_of_interest). Topic modeling is a word-level approach while most of the depression seeds in our lexicon are phrases. Consequently, seed phrase replacement plays an integral role in the success of our algorithm. Next, we apply platform-specific filtering, followed by non-ASCII character and stopword removal, as well as lemmatization. Platform-specific filtering includes substituting retweets (“RT @username” by RT), user mentions (“@username” by MENTION), and hyperlinks (by URL). For spelling correction, we utilize the PyEnchant spell checker library⁷. Furthermore, alphabetic character repetition (writing identical characters in sequence for emphasis, e.g., ftttttttt, slleeeeeep) is addressed by defining regular expressions and enhancing the available NLTK tweet tokenizer library⁸.

A. Qualitative Results

1) *Discovery of depressive symptoms:* Our ssToT model discovers depressive symptoms as latent topics from sliding window on buckets of timestamped tweets posted by users. We rank the top terms in each symptom $p(w|s)$ in descending order. Table II illustrates the sample of topics learned by ssToT and LDA model. The seeded words for the ssToT model are boldfaced, and words that are judged as relevant are italicized. We observe that by constraining seed terms to a specific symptom, the discovered terms are more relevant to that category. For example, in LDA model Topic 8 contains three terms relevant to “Sleep Disorder” (S3); however, it also contains lots of irrelevant terms which makes the emphasis of this topic off-target. Although Topic 6 from LDA contains terms relevant to “Eating Disorder” (S5), it also contains some terms related to “Sleep Disorder” and “Suicidal Thoughts” (S9). Similarly, for Topic 3, it contains terms associated with both the “Eating Disorder” and “Suicidal Thoughts” categories. Therefore, the topics discovered with LDA are not interpretable for the purpose of this study.

In contrast, the topics learned from the ssToT model contain more relevant terms associated with symptom category and more interpretable topics (see Table II). Additionally, the ssToT model also captures acronyms that people use in social media; for instance, in symptom 5 (Eating Disorder) “ugw” stands for “Ultimate Goal Weight” and “mfp” for “More Food Please”, or in symptom 2 (Lack of Interest) “idec” for “I Don’t Even Care”. We also observe the excessive usage of expressive interjections in language used by depressed users. Terms such

⁷<http://pythonhosted.org/pyenchant/>⁸<http://www.nltk.org/api/nltk.tokenize.html>

TABLE II. SAMPLE OF LEARNED TOPICS (SYMPTOMS) BY ssToT VS LDA AND THEIR ASSOCIATED COHERENCY SCORES

Model	Label	Top Words	UMass	UCI	NPMI
ssToT	Sleep Disorder	<i>cant sleep, wanna sleep, night, nighttime, sleepy, need to sleep, hour, sleepover, bedtime, go to sleep, mess, dream, midnight crying, painful, 5:00 AM, guilt, struggle, headaches tonight, morning, coffee, duvet, hungover, bbe</i>	-1.23	-1.28	-0.01
	Eating Disorder	<i>fat, eat, kg, weight loss, negative calories, lbs, thin, my thighs, paper thin, binge, eating disorder, abs, stomach, bulimic, hating, salad, pretend, gain, starve, mcdonalds, bones, chubby, flat, skip, wears, kcal, puffy, hippo, mfp, ugw</i>	-1.18	0.20	0.02
	Suicidal Thoughts	<i>self harm, cut, suicide, live, scar, blade, dead, alive, bleed, need my blade, death, hanging, deserve pain, kill me now, gun, want to die, knife, daisies, opinion, meh, razor, sharp, wrists, pictures, never wake up, wanna cut, stfu, ew</i>	-0.66	1.13	0.09
LDA	Topic 8	<i>Sleepover, september, lost, interest, exaggerating, its my fault, ugh, sleep, skin, dish, saved, wake up, blocked, blow, ipad, touches</i>	-1.44	-2.84	-0.1
	Topic 6	<i>thigh, blood, big, beautiful, thin, smile, sleep, blood, leave, stay, worthless, fat, tear, pretending, sadness, fake, ugly, god, skin, eat, morning</i>	-3.31	-2.65	-0.08
	Topic 3	<i>Blade, ugly, fat, blood, smile, mirror, call, fit, eat, stay, beautiful, sleep, big, tear, sad, devil, god, skin, music</i>	-2.69	-3.69	-0.09

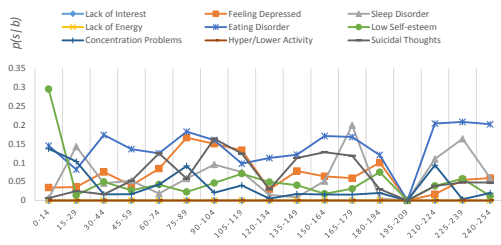


Fig. 1. Sample trend analysis of depressive symptoms

as “argh” (showing frustration), “aw” (indicative of disappointment), “feh” (indicative of feeling underwhelmed), “ew” (denoting disgust), “Huh” (indicator of confusion), “phew” (showing relief) were mostly discovered in their related symptoms category.

Furthermore, we observe that there are common themes and triggers of clinical depression at the community level that they do not exist in PHQ-9. In most cases, depressed users discuss their family and friend problems and the need for their support. For example the topic {family, hugs, attention, parents, competition, daddy, mums, sigh, grandma, losing, maam, friendless, love, friend, mommy, people, boyf, gf} shows that the person is suffering from a relationship problem. Another common theme is school and academic stress {schools, college, exam, classmate, friendless, teacher, assignment}.

To visualize depressive symptoms discovered for each user over a specific period, we keep the topics which contain at least certain number of seed terms as dominant words (among the top 20 terms associated with that topic) and discard the others.

Figure 1 depicts the sample of various depressive symptoms learned by ssToT and their distribution over time. It shows *when* and for *how long* a specific depressive symptom occurred. To check the validity of each topical trend, human annotators were asked to manually annotate all the buckets for the presence of all symptoms. Further details are provided in the next section.

VI. QUANTITATIVE RESULTS

Since ssToT is based on semi-supervised learning, it can be considered a clustering based approach and evaluated based on clustering evaluation measures. In addition, we are also able to employ classification accuracy to evaluate the performance of our method for symptom discovery. In this section, we first discuss the state of the art in clustering measures and how we adopt these measures to compare the performance of the ssToT model with existing methods. Then, we discuss the process of creating a ground truth dataset for evaluating the performance of the ssToT for discovering depressive symptoms as a multi-label classifier over different time periods.

Coherence Measures: Topic coherence measures score a single topic by identifying the degree of semantic similarity between high-scoring words in that topic. In this manner, we can distinguish between semantically interpretable topics and those which are artifacts of statistical inference. The state of the art for this evaluation criterion can be grouped into two major categories: intrinsic and extrinsic measures. Intrinsic measures evaluate the amount of information encoded by the topics over the original corpus used to train the topic models.

Another common intrinsic measure is UMass presented by [31], which measures the word co-occurrence in documents:

$$UMass(w_i, w_j) = \log \frac{D(w_i, w_j) + \epsilon}{D(w_i)}$$

where $D(w_i, w_j)$ counts the number of documents containing both w_i and w_j words, and $D(w_i)$ counts the ones containing w_i , over the same training corpus, and ϵ is the smoothing factor. The UMass metric computes these probabilities over the same corpus used to train the topic models.

Conversely, extrinsic evaluation metrics estimate the word co-occurrence statistics on external datasets such as Wikipedia. The UCI metric introduced by [32] utilizes the Pointwise Mutual Information (PMI) between two words,

$$UCI(w_i, w_j) = \log \frac{p(w_i, w_j) + \epsilon}{p(w_i)p(w_j)}$$

where the word probabilities are calculated by counting word co-occurrence in a sliding window over an external dataset such as Wikipedia. Recently, another topic coherence measurement has been introduced by [33] which considers

context vectors for every topic's top word. For every word w , a context vector is generated using word co-occurrence counts employing context window of size $+n$ surrounding that word. By calculating Normalized PMI (NPMI), they showed their method has a strong correlation with human topic coherence rating. The higher the topic coherence measure score, the higher the quality of the topics. This, in turn, leads to better topic interpretability, given that our purpose is to extract meaningful and interpretable topics associated with depressive symptoms. We used Palmetto⁹ for measuring the quality of topics learned based on NPMI and UCI measures (Wikipedia as an external corpus). UMass was measured by creating our Lucene index on tweets from users in set U (intrinsic evaluation)¹⁰. Table II shows a sample of the coherency of topics learned (symptoms) for LDA and ssToT models based on UMass, UCI, and NPMI metrics. We can clearly see that the topics learned by the ssToT model are more coherent for all the three measures.

Baselines: To further evaluate the ssToT-learned topics, we compare them with the topics obtained from a set of existing unsupervised and semi-supervised approaches. **k-means:** A clustering approach based on distributional similarity employing cosine similarity measure. **LSA:** An unsupervised approach that gleams distributional semantics by clustering correlated terms into latent topics using singular value decomposition. **LDA:** A Bayesian approach that represents a document as a mixture of topics. **BTM:** A state-of-the-art unsupervised topic modeling framework for short texts which utilizes distributed representations of words and phrases [34]. **Partially Labeled LDA:** A semi-supervised topic model which constrains latent topics to align them with human-provided labels [24].

To determine the number of topics for all LDA variants, we use perplexity using 80% of the data to train and 20% to test. We choose 15 topics as a proper level of granularity as it has the lowest perplexity and is suitable for our task. We set the number of Gibbs iterations to 1,000, α to 0.5, β to 0.1 and the rest of the parameters to default values. We use the Stanford Topic Modeling Toolbox¹¹ to run all LDA variants except BTM, which is downloaded from its author's webpage¹². We use cosine similarity as a distance function for k-means. Table III denotes the average coherence score for each model. Due to space limitations, we only report the average coherence of all symptoms for each algorithm. Coherency measures judge each model's output based on how well they represent a specific topic. This aligns with our objective of providing outputs that are well associated with depressive symptoms rather than some generic set of terms grouped together. These numbers indicate that the ssToT model outperforms other state-of-the-art techniques regardless of the corpus that probabilities are gained from: Wikipedia for UCI or the same corpus in UMass. We note that although, on average, the ssToT model outperforms the other five baselines in terms of discovering coherent topics, there are rare exceptions. For instance, the topic containing {fat, time, feel, dinner, weight, eat, hate, skinny} learned by the BTM algorithm about Eating Disorder has scores of -0.19, 0.8, and 0.08 for UMass, UCI, and NPMI respectively, which

⁹<http://rebrand.ly/palme9bf7>

¹⁰<http://rebrand.ly/howtod52e>

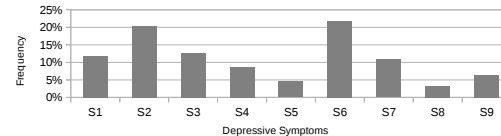
¹¹<http://nlp.stanford.edu/software/tmt/tmt-0.4/>

¹²<https://github.com/xiaohuiyan/BTM>

TABLE III. AVERAGE COHERENCY OF DIFFERENT MODELS VS. ssToT

Model	UMass	UCI	NPMI
LDA	-2.68	-3.03	-0.109
BTM	-1.42	-2.18	-0.058
P-LDA	-1.12	-2.48	-0.123
K-Mean	-1.70	-2.95	-0.102
LSA	-1.43	-3.23	-0.107
ssToT	-1.00	-1.62	-0.026

Fig. 2. Symptom distribution in the gold standard dataset



implies that it is more coherent compared to its associated topics learned by ssToT (see Table III). However, when we further analyzed the rest of the topics learned by BTM, we noticed the poor quality of the other learned topics, as the average coherency also indicates.

VII. SYMPTOM PREDICTION (MULTI-LABEL CLASSIFICATION)

We showed how ssToT is highly effective in terms of the quality of learned topics (depressive symptoms). In this section, we further investigate the power of the ssToT model as a multi-label classifier. Specifically, we try to predict the correct set of labels (depressive symptoms) for each bucket of tweets. We build a ground truth dataset of 10400 tweets in 192 buckets. Each bucket contains tweets that are posted by the user within span of 14 days (in compliance with PHQ-9). Tweets are selected from a randomly sampled subset of both self-reported depressed users and random users. Three human judges (undergraduate students who are native English speakers) manually annotated each tweet using the nine PHQ-9 categories as labels. Additionally, the non-relevant tweets that do not show any depressive symptoms, have been labeled as “cannot judge”. The average inter-annotator agreement is $K=0.74$ based on Cohen's Kappa statistics. We build a labeled bucket by merging the labeled tweets. Figure 2 depicts symptoms distribution in the gold standard dataset.

Our semi-supervised ssToT model does not use the labeled data during training. In particular, we are not supervising the LDA model with a labeled-dataset like Labeled-LDA [35]. Instead, we are using the labeled dataset for evaluating the performance of the ssToT model in assigning a set of symptoms to each bucket. We evaluate the performance of the ssToT in terms of the average precision, recall, and F-score in detecting depressive symptoms for each bucket of tweets when tested against human judgment. We show that our ssToT model can predict the *presence* of each of the nine depressive symptoms for a different bucket with an accuracy of 0.68 (see Table IV) which gives a precision of 0.72 *on average*.

We observed that the best results were achieved for “Lack of Interest” (symptom 1) with a 0.90 F-Measure and the worst result was obtained for “Concentration Problems” (symptom 7). We noticed that our ssToT model works well with less descriptive symptoms since it generates relevant seed terms to discover the latent symptom. For instance, consider the

following tweets selected from our dataset: “Overthinking always destroy my mood”, “This essay is dragging so much, can’t deal with essays and revisions any more :(”, “I need a break from my thoughts”, and “my head is such a mess right now”. The first tweet contains the depression-indicative keywords “overthinking”, which can easily be interpreted as “Concentration Problems” while the other utterances, although they are labeled into the same symptom category (according to our human annotator), are descriptive and metaphoric and do not contain any depressive-indicative term. However, there are some tweets such as “overthinking killed my happiness” that, even though they contain the depression-indicative terms of “Concentration Problems”, (overthinking in this case) cannot be grouped into “Concentration Problems” category. Such examples contribute to a high number of false positives and low precision for this category. Furthermore, sometimes correctly determining the category of depressive symptoms is challenging even for a human. For instance, in the tweet “Need to sleep, always so f***ing tired” one may categorize it as “Lack of Energy” while another may consider it as “Sleep Disorder”.

To further test the robustness of our ssToT model as a multi-label classifier, we compare its results to common supervised approaches for performing multi-label classification, namely the *binary relevance* (BR) and *classifier chains* (CC) methods [36]. The BR method transforms the problem of multi-label into multiple binary models by creating one model for each label. For this, each binary model will be trained to predict the relevance of each of the labels. On the other hand, CC is a chaining method that uses L binary transformations (one for each label) similar to BR, but it can also model label correlations while maintaining acceptable computational complexity. As supervised baseline approaches, Multinomial Naive Bayes and SVM models have been chosen for the two aforementioned methods. These two models have been widely utilized as a baseline by most previous studies [24]. Note that in the task of supervised multi-label classification, labels are available during training. We used Meka (a Multi-label Extension to WEKA)¹³ for building the baselines. We use a bag-of-words model and perform 10-fold cross-validation to evaluate accuracy for each symptom (see Table IV).

These results show that in spite of the semi-supervised nature of ssToT model, it is competitive with supervised approaches and improves upon them in five out of nine symptom classification in terms of F-score along with providing better averaged accuracy. A key advantage of ssToT over supervised approach is that it does not require labor-intensive, expensive, and time-consuming manual annotation of data in training.

Our study has limitations. For users who do not generate ample content on their profiles or are reluctant to publicly reveal their depressive symptoms, we cannot assess their depressive behaviors. Additionally, we only detect the presence, duration, and frequency of symptoms rather than their severity. Furthermore, more severely depressed individuals may be more inclined to publicly express their depression and biasing our sample.

¹³<http://meka.sourceforge.net/>

TABLE IV. MODEL’S PERFORMANCE FOR BUCKET LEVEL SYMPTOM PREDICTION, (P:PRECISION, R:RECALL, F:F-SCORE, AA: IS THE AVERAGE ACCURACY FOR EACH MODEL.

Model	AA.		S1	S2	S3	S4	S5	S6	S7	S8	S9
BR-MNB	0.66	P	0.68	0.94	0.67	0.62	0.71	0.96	0.63	0.38	0.96
		R	0.74	0.78	0.72	0.59	0.86	0.95	0.82	0.86	0.93
		F	0.74	0.81	0.70	0.60	0.89	0.90	0.73	0.84	0.96
CC-MNB	0.63	P	0.73	0.95	0.73	0.61	0.90	0.99	0.70	0.71	0.91
		R	0.66	0.80	0.68	0.55	0.84	0.87	0.65	0.85	0.80
		F	0.72	0.83	0.73	0.57	0.91	0.92	0.72	0.90	0.88
BR-SVM	0.695	P	0.88	0.81	0.91	0.94	0.98	0.74	0.62	0.81	0.79
		R	0.69	0.94	0.69	0.42	0.84	0.96	0.80	0.80	0.82
		F	0.79	0.85	0.76	0.53	0.91	0.93	0.80	0.88	0.90
CC-SVM	0.71	P	0.87	0.93	0.79	0.93	0.97	0.98	0.74	0.92	1.0
		R	0.70	0.93	0.66	0.40	0.83	0.94	0.79	0.80	0.84
		F	0.80	0.86	0.75	0.51	0.91	0.95	0.82	0.89	0.91
ssToT	0.68	P	0.87	0.93	0.79	0.91	0.97	0.98	0.74	0.92	0.79
		R	0.93	0.98	0.69	0.53	0.90	0.92	0.19	0.24	0.59
		F	0.90	0.89	0.78	0.68	0.93	0.82	0.30	0.38	0.68

VIII. CONCLUSION AND FUTURE WORK

We demonstrated the impact of social media on extraction and timely monitoring of depression symptoms. We developed a statistical model using a hybrid approach that combines a lexicon-based technique with a semi-supervised topic modeling technique to extract per user topic distribution (clinical, symptomatic of depression) and per topic word distribution (symptom indicators) by textual analysis of tweets over different time windows. Our approach complements the current questionnaire-driven diagnostic tools by gleaned depression symptoms in a continuous and unobtrusive manner. Our experimental results reveal that there are significant differences in the topic preferences and word usage pattern of the self-declared depressed group from random users in our dataset which indicates the competency of our model for this task. Our model yields promising results with an accuracy of 68% and a precision of 72% for capturing depression symptoms per user over a time interval which is competitive with a fully supervised approach. In future, we plan to apply our approach to various data sources such as longitudinal electronic health record (EHR) systems and private insurance reimbursement and claims data, to develop a robust “big data” platform for detecting clinical depressive behavior at the community level.

IX. ACKNOWLEDGEMENT

We are thankful to Surendra Marupudi and Ankita Saxena for helping us with data collection. We also thank Jibril Ikhara for his proofreading. Research reported in this publication was supported in part by NIMH of the National Institutes of Health (NIH) under award number R01MH105384-01A1.

REFERENCES

- [1] M. G. Haselton, D. Nettle, and D. R. Murray, “The evolution of cognitive bias,” *The handbook of evolutionary psychology*, 2005.
- [2] 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 *Proceedings of the 2nd Workshop on Computational Linguistics and Clinical Psychology: From Linguistic Signal to Clinical Reality*, 2015.
- [3] M. Mitchell, K. Hollingshead, and G. Coppersmith, “Quantifying the language of schizophrenia in social

- media,” in *Proceedings of the 2015 Annual Conference of the North American Chapter of the ACL: Human Language Technologies (NAACL HLT)*, 2015.
- [4] A. Benton, M. Mitchell, and D. Hovy, “Multitask learning for mental health conditions with limited social media data,” *EACL*, 2017.
 - [5] G. C. M. D. C. Harman, “Quantifying mental health signals in twitter,” *ACL 2014*, 2014.
 - [6] A. H. Yazdavar, H. S. Al-Olimat, T. Banerjee, K. Thirunarayan, and A. P. Sheth, “Analyzing clinical depressive symptoms in twitter,” 2016.
 - [7] Y. Neuman, Y. Cohen, D. Assaf, and G. Kedma, “Proactive screening for depression through metaphorical and automatic text analysis,” *Artificial intelligence in medicine*, vol. 56, no. 1, pp. 19–25, 2012.
 - [8] C. Karmen, R. C. Hsiung, and T. Wetter, “Screening internet forum participants for depression symptoms by assembling and enhancing multiple nlp methods,” *Computer methods and programs in biomedicine*, 2015.
 - [9] M. De Choudhury, M. Gamon, S. Counts, and E. Horvitz, “Predicting depression via social media,” in *ICWSM*.
 - [10] T. Nguyen, D. Phung, B. Dao, S. Venkatesh, and M. Berk, “Affective and content analysis of online depression communities,” *IEEE Transactions on Affective Computing*, vol. 5, no. 3, pp. 217–226, 2014.
 - [11] G. Coppersmith, M. Dredze, C. Harman, K. Hollingshead, and M. Mitchell, “Clpsych 2015 shared task: Depression and ptsd on twitter,” in *Proceedings of the 2nd Workshop on Computational Linguistics and Clinical Psychology*, 2015.
 - [12] 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,” *NAACL HLT 2015*, 2015.
 - [13] K. Kroenke, R. L. Spitzer, and J. B. Williams, “The phq-9,” *Journal of general internal medicine*.
 - [14] M. De Choudhury, S. Counts, and E. Horvitz, “Social media as a measurement tool of depression in populations,” in *Proceedings of the 5th Annual ACM Web Science Conference*. ACM, 2013, pp. 47–56.
 - [15] X. Wang, C. Zhang, Y. Ji, L. Sun, L. Wu, and Z. Bao, “A depression detection model based on sentiment analysis in micro-blog social network,” in *Pacific-Asia Conference on KD and DM*. Springer, 2013, pp. 201–213.
 - [16] H.-H. Shuai, C.-Y. Shen, D.-N. Yang, Y.-F. Lan, W.-C. Lee, P. S. Yu, and M.-S. Chen, “Mining online social data for detecting social network mental disorders,” in *Proceedings of the 25th International Conference on World Wide Web*, 2016, pp. 275–285.
 - [17] D. Preotiuc-Pietro, J. Eichstaedt, G. Park, M. Sap, L. Smith, V. Tobolsky, H. A. Schwartz, and L. Ungar, “The role of personality, age and gender in tweeting about mental illnesses,” in *NAACL HLT*, 2015.
 - [18] P. Thompson, C. Poulin, and C. J. Bryan, “Predicting military and veteran suicide risk: Cultural aspects,” in *Proceedings of the Workshop on Computational Linguistics and Clinical Psychology*, 2014.
 - [19] G. Coppersmith, K. Ngo, R. Leary, and A. Wood, “Exploratory analysis of social media prior to a suicide attempt,” in *Proceedings of the Third Workshop on Computational Linguistics and Clinical Psychology*, 2016.
 - [20] D. N. Milne, G. Pink, B. Hachey, and R. A. Calvo, “Clpsych 2016 shared task: Triaging content in online peer-support forums,” in *Proceedings of the Third Workshop on Computational Linguistics*, 2016.
 - [21] M. Park, D. W. McDonald, and M. Cha, “Perception differences between the depressed and non-depressed users in twitter,” in *ICWSM*, 2013.
 - [22] T. Wang, Z. Huang, and C. Gan, “On mining latent topics from healthcare chat logs,” *Journal of Bio Info*, 2016.
 - [23] D. Andrzejewski and X. Zhu, “Latent dirichlet allocation with topic-in-set knowledge,” in *Proceedings of the NAACL HLT 2009 Workshop on Semi-Supervised Learning for Natural Language Processing*. ACL, 2009.
 - [24] D. Ramage, C. D. Manning, and S. Dumais, “Partially labeled topic models for interpretable text mining,” in *Proceedings of the 17th ACM SIGKDD international conference on KD and DM*. ACM, 2011, pp. 457–465.
 - [25] D. Andrzejewski, X. Zhu, M. Craven, and B. Recht, “A framework for incorporating general domain knowledge into latent dirichlet allocation using first-order logic,” in *IJCAI*, no. 1, 2011.
 - [26] D. Andrzejewski, X. Zhu, and M. Craven, “Incorporating domain knowledge into topic modeling via dirichlet forest priors,” in *Proceedings of the 26th Annual International Conference on Machine Learning*. ACM, 2009.
 - [27] A. Mukherjee and B. Liu, “Aspect extraction through semi-supervised modeling,” in *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics: Long Papers-Volume 1*. Association for Computational Linguistics, 2012, pp. 339–348.
 - [28] A. H. Yazdavar, M. Ebrahimi, and N. Salim, “Fuzzy based implicit sentiment analysis on quantitative sentences,” *arXiv preprint arXiv:1701.00798*, 2017.
 - [29] M. Ebrahimi, M. Ebrahimi, A. H. Yazdavar, A. H. Yazdavar, N. Salim, N. Salim, S. Eltyeb, and S. Eltyeb, “Recognition of side effects as implicit-opinion words in drug reviews,” *Online Information Review*, vol. 40, no. 7, pp. 1018–1032, 2016.
 - [30] M. Ebrahimi, A. H. Yazdavar, and A. Sheth, “On the challenges of sentiment analysis for dynamic events.”
 - [31] D. Mimno, H. M. Wallach, E. Talley, and A. McCallum, “Optimizing semantic coherence in topic models,” in *Empirical Methods in NLP*. ACL, 2011.
 - [32] D. Newman, J. H. Lau, K. Grieser, and T. Baldwin, “Automatic evaluation of topic coherence,” in *Human Language Technologies: The 2010 Conference of the North American Chapter of the ACL*, 2010.
 - [33] N. Aletas and M. Stevenson, “Evaluating topic coherence using distributional semantics,” in *Proceedings of the 10th International Conference on Computational Semantics (IWCS 2013)-Long Papers*, 2013.
 - [34] X. Yan, J. Guo, Y. Lan, and X. Cheng, “A bitern topic model for short texts,” in *Proceedings of the 22nd international conference on World Wide Web*. ACM, 2013, pp. 1445–1456.
 - [35] D. Ramage, D. Hall, R. Nallapati, and C. D. Manning, “Labeled lda: A supervised topic model for credit attribution in multi-labeled corpora,” in *Proceedings of the 2009 Conference on Empirical Methods in NLP: Volume 1-Volume 1*. ACL, 2009, pp. 248–256.
 - [36] J. Read, B. Pfahringer, G. Holmes, and E. Frank, “Classifier chains for multi-label classification,” *Machine learning*, 2011.

Mental Health Analysis via Social Media Data

Amir Hossein Yazdavar
Kno.e.sis Center
 Wright State University
 Dayton, OH, USA
 amir@knoesis.org

Mohammad Saied Mahdavejad
Kno.e.sis Center
 Wright State University
 Dayton, OH, USA
 msmn.nova@gmail.com

Goonmeet Bajaj
Kno.e.sis Center
 Wright State University
 Dayton, OH, USA
 goonmeet@knoesis.org

Krishnaprasad Thirunarayan
Kno.e.sis Center
 Wright State University
 Dayton, OH, USA
 tkprasad@knoesis.org

Jyotishman Pathak
Division of Health Informatics
 Cornell University
 New York, NY, USA
 jyp2001@med.cornell.edu

Amit Sheth
Kno.e.sis Center
 Wright State University
 Dayton, OH, USA
 amit@knoesis.org

Abstract—With ubiquity of social media platforms, millions of people are routinely sharing their moods, feelings and even their daily struggles with mental health issues by expressing it verbally or indirectly through images they post. In this study, we aim to examine exploitation of big multi-modal social media data for studying depressive behavior and its population trend across the U.S. to better understand a regions influence on the prevailing environment and available care. In particular, employing statistical techniques along with the fusion of heterogeneous features gleaned from different modalities (shared images and textual content), we build models to detect depressed individuals and their demographics.

Keywords—Multi-modal Analysis, Machine Learning, Natural Language Processing, Statistical analysis, Social Media, Mental Health, Regression

I. INTRODUCTION

Depression is a highly prevalent public health challenge and a major cause of disability worldwide. According to the World Mental Health Survey conducted in 17 countries, on average, about 5% of people reported having an episode of depression in 2011[1]. It also affects 6.7% of Americans (that is, more than 16 million) each year¹. Untreated or under-treated clinical depression can be serious enough to lead to suicide and other risky behaviors such as drug or alcohol addiction². More than 90% of people who commit suicide have been diagnosed with depression [2].

A global effort to curb clinical depression involves identifying it through survey-based methods via phone or online questionnaires. These approaches mainly suffer from under-representation as well as sampling biases (a small group of respondents). Besides, survey data also exhibit problems due to temporal gaps between the data collection and dissemination of findings, reflecting participant's responses over

a short period of time. In contrast, with an unprecedented growth of social media, large number of people voluntarily share large amounts of data by expressing their moods, feelings, emotions, and daily struggles with mental health problems on social media platforms like Twitter. This offers opportunities for new understanding of these communities. For instance, the news headlines such as "Twitter Fail: Teen Sent 144 Tweets Before Committing Suicide & No One Helped" highlights the need for better tools for gleaning useful insights from user generated content on social media.

Recent years have witnessed rapid growth in the analysis of social media for studying a wide range of health problems, from detecting influenza epidemic [3], [4] and cardiac arrest [5] to study mood and mental health conditions [6], [7]. Previous research efforts suggested that language style, sentiment, ego-network, and user engagement are discriminating features to predict the likelihood of depression in a post [8] or in an individual [9], [10]. These studies often represent psychological status of online users via their language via psycholinguistic analysis, supervised and unsupervised language modeling, or studying individuals topic of interest. However, except few attempts [cite, cite, cite] investigations in the fields are seldom concerned with visual attributes of mental health reflects on various social media platforms. Interestingly, according to eMarketer, photos accounted for 75% of content posted by Facebook pages worldwide and they are the most engaging type of content on Facebook (87%). It has been argued that sharing them is the best way get more attention from the followers. Indeed, an old saying of "a picture is worth a thousand words" recently changed to "photos are worth a million likes". Similar to Facebook, photos are more engaging for Twitter users. The tweets attached with image links get two times more engagement rate of those without 4. We recall that getting social support from peers has been highlighted as a primary motivation for sharing depressive indicative content in social media [6].

¹<http://bit.ly/2okBKNy>

²<http://www.webmd.com/depression/guide/untreated-depression-effects#1>

The easiness of expressing emotion through images where they often gaining more attention compared to the verbal form, is plausible motivation for sharing depressive images. Besides, as the psychologist Carl Rogers highlights, we often pursue attitudes which bring us closer to our Ideal-Self. In this regard, the choice of profile image can either represent our online persona or persona we choose to paint for others to see. We believe this can have roots in the mental health status of a person, and the visual attributes of it can provide emotional expression that can yield insights into mental illness. Inspired by that, we study the profile pictures of likely depressed individuals in order to discover relevant signals from colors, aesthetic and facial attributes, to better understand the psychology and the intent behind choosing the personal profile image.

Furthermore, the recent advancements in deep-convolutional neural networks, specifically for the image analysis task, has lead to a significant improvement in age and gender classification. Inspired by that, we develop a big data approach to automatically detect likely depressed individuals by exploiting heterogeneous set of features gleaned from different modality from studying the content of posted images (colorfulness, hue variance, sharpness, brightness, blurriness, naturalness), the choice of profile picture (gender, age, and emotion estimation), the choice of screen name, language features from both generated textual content and profiles description (n-gram, emotion, sentiment [11], [12], [13])(see Figure 1). In particular, we address the following research questions: 1) How well do the content of posted images (colors, aesthetic and facial presentation) can reflect depressive behavior? Does the choice of profile picture show any psychological trait of online depressed persona? Are they reliable enough to represent the demographics such as age and gender? 2) Are there any underlying common themes among depressed individuals generated visual and textual content?

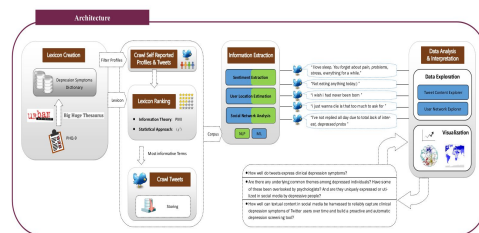


Figure 1. System architecture and the overall process

ACKNOWLEDGMENT

Research reported in this publication was supported in part by NIMH of the National Institutes of Health (NIH)

under award number R01MH105384-01A1.

REFERENCES

- [1] M. Marcus, M. T. Yasamy, M. van Ommeren, D. Chisholm, S. Saxena *et al.*, "Depression: A global public health concern," *WHO Department of Mental Health and Substance Abuse*, vol. 1, pp. 6–8, 2012.
- [2] M. D. Rudd, A. L. Berman, T. E. Joiner Jr, M. K. Nock, M. M. Silverman, M. Mandrusiak, K. Van Orden, and T. Witte, "Warning signs for suicide: Theory, research, and clinical applications," *Suicide and Life-Threatening Behavior*, vol. 36, no. 3, pp. 255–262, 2006.
- [3] A. Culotta, "Lightweight methods to estimate influenza rates and alcohol sales volume from twitter messages," *Language resources and evaluation*, vol. 47, no. 1, pp. 217–238, 2013.
- [4] E. Aramaki, S. Maskawa, and M. Morita, "Twitter catches the flu: detecting influenza epidemics using twitter," in *Proceedings of the conference on empirical methods in natural language processing*. Association for Computational Linguistics, 2011, pp. 1568–1576.
- [5] J. C. Bosley, N. W. Zhao, S. Hill, F. S. Shofer, D. A. Asch, L. B. Becker, and R. M. Merchant, "Decoding twitter: Surveillance and trends for cardiac arrest and resuscitation communication," *Resuscitation*, vol. 84, no. 2, pp. 206–212, 2013.
- [6] G. Coppersmith, M. Dredze, C. Harman, K. Hollingshead, and M. Mitchell, "Clpsych 2015 shared task: Depression and ptsd on twitter," in *Proceedings of the 2nd Workshop on Computational Linguistics and Clinical Psychology: From Linguistic Signal to Clinical Reality*, 2015, pp. 31–39.
- [7] G. Coppersmith, C. Harman, and M. Dredze, "Measuring post traumatic stress disorder in twitter," in *ICWSM*, 2014.
- [8] M. De Choudhury, M. Gamon, S. Counts, and E. Horvitz, "Predicting depression via social media," in *ICWSM*, 2013, p. 2.
- [9] A. H. Yazdavar, H. S. Al-Olimat, M. Ebrahimi, G. Bajaj, T. Banerjee, K. Thirunarayan, J. Pathak, and A. Sheth, "Semi-supervised approach to monitoring clinical depressive symptoms in social media," in *Proceedings of the 2017 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining 2017*. ACM, 2017, pp. 1191–1198.
- [10] A. H. Yazdavar, H. S. Al-Olimat, T. Banerjee, K. Thirunarayan, and A. P. Sheth, "Analyzing clinical depressive symptoms in twitter," 2016.
- [11] M. Ebrahimi, A. H. Yazdavar, N. Salim, and S. Eltyeb, "Recognition of side effects as implicit-opinion words in drug reviews," *Online Information Review*, vol. 40, no. 7, pp. 1018–1032, 2016.
- [12] A. H. Yazdavar, M. Ebrahimi, and N. Salim, "Fuzzy based implicit sentiment analysis on quantitative sentences," *arXiv preprint arXiv:1701.00798*, 2017.
- [13] M. Ebrahimi, A. H. Yazdavar, and A. Sheth, "Challenges of sentiment analysis for dynamic events," *IEEE Intelligent Systems*, vol. 32, no. 5, pp. 70–75, 2017.