Then and Now: Quantifying the Longitudinal Validity of Self-Disclosed Depression Diagnoses

Keith Harrigian and Mark Dredze

Johns Hopkins University

kharrigian@jhu.edu, mdredze@cs.jhu.edu

Abstract

Self-disclosed mental health diagnoses, which serve as ground truth annotations of mental health status in the absence of clinical measures, underpin the conclusions behind most computational studies of mental health language from the last decade. However, psychiatric conditions are dynamic; a prior depression diagnosis may no longer be indicative of an individual's mental health, either due to treatment or other mitigating factors. We ask: to what extent are self-disclosures of mental health diagnoses actually relevant over time? We analyze recent activity from individuals who disclosed a depression diagnosis on social media over five years ago and, in turn, acquire a new understanding of how presentations of mental health status on social media manifest longitudinally. We also provide expanded evidence for the presence of personality-related biases in datasets curated using self-disclosed diagnoses. Our findings motivate three practical recommendations for improving mental health datasets curated using self-disclosed diagnoses:

- 1. Annotate diagnosis dates and psychiatric comorbidities
- 2. Sample control groups using propensity score matching
- 3. Identify and remove spurious correlations introduced by selection bias

1 Introduction

The ability to provide equitable access to psychiatric healthcare has become more difficult than ever, inhibited by an entanglement of lingering public policy effects (Miranda et al., 2020), heightened levels of physician burnout (Johnson et al., 2018), and infrastructural challenges arising from global crisis (Davis et al., 2021). Meanwhile, social media platforms have become the predominant means of communication for much of the population, providing the opportunity to share personal experiences

and seek support from others (Mueller et al., 2021). Noting these parallel timelines, computational scientists have devoted substantial effort to engineering statistical models capable of translating social media data into reliable insights regarding mental health. Core objectives of this work include optimizing psychiatric treatment, identifying early stages of mental illness, and measuring the effect of public policy on a population's well-being (Losada et al., 2017; Fine et al., 2020).

The most significant advances in computational mental health research have not come from improved modeling architectures (Benton et al., 2017b), but from methods for curating largescale datasets which contain robust and clinicallyrelevant ground truth annotations of mental health status (Coppersmith et al., 2014). Use of regular expressions to identify genuine self-disclosures of a psychiatric diagnosis remains one of the most widely adopted annotation mechanisms by the research community (Chancellor and De Choudhury, 2020; Harrigian et al., 2021), offering a relatively reliable proxy in place of clinical measures which are not only costly to collect, but also often unable to be shared beyond a single institution due to patient privacy policies (Macavaney et al., 2021). Datasets leveraging self-disclosed diagnoses as annotations of mental health status have yielded a variety of insights that align with clinical knowledge and psychological theory (Mowery et al., 2017; Lee et al., 2021). However, a growing body of work has raised questions about whether such datasets provide sufficient information to train statistical models that generalize to new populations (Harrigian et al., 2020; Aguirre et al., 2021).

Despite the prevalence of datasets dependent on self-disclosure, no analyses have considered how associating a single self-disclosed diagnosis label with data from a variable-length period of time may inhibit the learning of robust statistical relationships. If a user tweets a depression diagnosis in 2015, is their data from 2018 still representative of the condition? Presentation of several mental health conditions change dynamically and (sometimes) precipitously over time (Collishaw et al., 2004). Yet, it remains common in the computational research community to treat mental health conditions as a static attribute with equal relevance at multiple time points (MacAvaney et al., 2018). In reality, it is likely that only a small fraction of an individual's social media activity is appropriate for training optimal classifiers. Moreover, that a mental health status label may be appropriate for only a subset of time suggests that evaluations of longitudinal model generalization as they are traditionally structured in the community may be insufficient (Sadeque et al., 2018).

We ask: to what extent do mental health diagnosis self-disclosures remain valid over time? We focus specifically on extended durations (i.e., multiple years), a setting which has particular relevance to those who wish to estimate generalization strength of their statistical classifiers for use in longitudinal monitoring applications, as well as those interested in updating existing models with new data to mitigate the effects covariate shift (Agarwal and Nenkova, 2021). In reviewing recent online activity from individuals in the 2015 CLPsych Shared Task dataset who disclosed a depression diagnosis on Twitter over five years ago (Coppersmith et al., 2015), we not only acquire a new understanding of how presentations of mental health status on social media present over time, but also find new evidence to support prior claims regarding the presence of personality-related confounds in datasets curated using self-disclosures (Preotiuc-Pietro et al., 2015; Vukojevic and Šnajder, 2021). Our analysis provides critical guidance to practitioners as they curate mental health social media datasets, while also elucidating factors which inhibit robustness in a dataset that remains one of the most widely adopted by the research community.

2 Background

The majority of mental health research based on social media leverages the same experimental design—assume individuals have a fixed mental health status and attempt to infer this latent attribute using historical online activity traces (e.g., posts, follower network dynamics) (Guntuku et al., 2017; Chancellor and De Choudhury, 2020). This training setting is convenient given the inherent

Dataset	Dates	# Users	# Posts
Original	2012 – 2015	D: 477 C: 872	D: 1,121,388 C: 1,907,508
Updated	2012 - 2021	D: 444 C: 172	D: 1,372,868 C: 546,826

Table 1: Summary statistics for the original and updated versions of the 2015 CLPsych Shared Task dataset, further stratified by [C]ontrol and [D]epression groups.

complexities of acquiring temporally-granular psychiatric measures at scale (Canzian and Musolesi, 2015). However, the setting implicitly relies on assumptions that are not supported by clinical knowledge regarding psychiatric dynamics (Johnson and Nowak, 2002; Schoevers et al., 2005). Some work has been done to incorporate time-based priors into mental health models, which allow practitioners to train statistical classifiers using a static label while also explicitly accounting for longitudinal variation in label relevance (Wongkoblap et al., 2019; Uban et al., 2021). Others have eschewed the use of a static label altogether and instead curated datasets that contain multiple points of ground truth mental health status, albeit still with some element of historical data aggregation (Chancellor et al., 2016).

Temporally-aware classifiers have achieved better performance benchmarks than their static counterparts in some cases (Rao et al., 2020), though these evaluations remain limited by the dearth of data with mental health status annotations at multiple time points. Meanwhile, datasets which do support dynamic evaluation are curated almost exclusively using protected clinical measures (Reece et al., 2017), cost-intensive interviews (Nobles et al., 2018), or non-trivial shifts in non-language-based online behavior (De Choudhury et al., 2016).

Computational studies that have focused on self-disclosed diagnoses have not comprehensively reviewed how individual activity evolves over long periods of time (Saha et al., 2021). Our study thus fulfills an important void in the research space by providing a new understanding of long term mental health dynamics in social media, and more particularly, within convenience samples curated using self-disclosed diagnoses.

3 Data

We support our study using a newly updated version of the 2015 CLPsych Shared Task dataset

(Coppersmith et al., 2015). The original Twitter dataset was constructed in a two-stage process, with regular-expressions first being used to identify candidate self-disclosures of a depression diagnosis and experts manually verifying the authenticity of the match thereafter. Individuals in the control group were sampled randomly from the 1% public Twitter stream such that the joint distribution of inferred age and gender attributes (Sap et al., 2014) was in alignment with the depression group. Up to 3,000 tweets were acquired for each individual in the resulting sample using Twitter's public API. The dataset has not only become one of the most widely adopted social media datasets for mental health (Harrigian et al., 2021), but also inspired the annotation procedures for numerous successors across various platforms and languages (Cohan et al., 2018; Shen et al., 2018).

In line with guidance from Benton et al. (2017a), individual identifiers in the official version of the CLPsych dataset have been anonymized, with linkages between anonymized and de-anonymized identifiers erased in entirety. However, the original de-anonymized identifiers remain available under explicit permission from Coppersmith et al. (2015), who provided this information to reverse engineer the original anonymization mapping. To do so, we first query up to 3,200 of the most recent tweets from each de-anonymized user identifier using Twitter's public API and further isolate all relevant tweets found in our institution's cache of Twitter's 1% data stream. We identify candidate pairs of anonymized and de-anonymized accounts based on overlap of raw timestamps within the original dataset's collection window. Normalized text (i.e., punctuation removal, case standardization) from candidate pairs is compared using exact matching to verify final linkages.

Statistics for the original dataset and its updated counterpart are provided in Table 1. We find that a majority of accounts which were unable to be linked had significantly smaller activity traces in the original dataset. These accounts are likely to either have been deleted in entirety or to have tweeted with a small enough frequency such that the 1% stream does not contain any samples. The discrepancy in match rates between individuals in the depression and control groups is unfortunately not fully-understood, though discussions with the dataset's authors suggest this may just be an artifact of the original archival process.

Preprocessing. Twitter's language tags and automatic language identification (Lui and Baldwin, 2012) are used to isolate English text. Retweets are excluded to most acutely highlight personal experiences with depression over time. Unless specified otherwise, keyword-based tweet filtering is applied to preemptively mitigate sampling-induced biases which can artificially inflate estimates of predictive performance. Some of these biases have been recognized and addressed by the research community (e.g., filtering tweets which include diagnosis disclosures and/or mental health related keywords/hashtags) (De Choudhury and De, 2014), while others have been traditionally overlooked.

A preliminary qualitative analysis of influential n-grams and their source tweets reveals a previously unrecognized surplus of "fan accounts" (e.g., supporters of Harry Styles and Demi Lovato) and tweets containing account statistics (e.g., new followers) within the depression cohort. Meanwhile, daily horoscope tweets were identified with an anomalous frequency within the control group. The latter two sources of noise do not have a clear clinical explanation, while the former (i.e., fan accounts) arises in the context of discussion regarding the mental health of young celebrities. Although some of these motifs represent genuine behavioral correlates of depression, their importance in prediction tends to be inflated due to context of the original collection time period.

4 Inference Under Latent Dynamics

Enabling reliable use of statistical models to evaluate change in mental health status remains a core objective for computational researchers (Choi et al., 2020; Fine et al., 2020). Our success in this task domain critically depends on access to ground truth at multiple time points, not only for evaluating generalization error (DeMasi et al., 2017; Tsakalidis et al., 2018), but also for mitigating the effects of covariate shift (Sugiyama and Kawanabe, 2012). As discussed above, it is often trivial to update activity traces for individuals with a prior mental health diagnosis disclosure. Nonetheless, clinical knowledge suggests original disclosure-based labels may not be relevant over the course of time, either due to a condition's episodic presentations (Angst et al., 2009) or the effects of psychiatric treatment (Saha et al., 2021). We ask whether the CLPsych Shared Task dataset supports this theory.

Methods. A natural framework for answering

this inquiry emerges from computational research regarding label noise (Frénay and Verleysen, 2013). Under such a perspective, we can view changes in mental health status as a stochastic process which blindly alters the correctness of class labels over time. The implications of this mechanism allow us to reason about predictive performance of a statistical classifier within and outside of the time period in which it is trained. Differences in within-timeperiod performance for two different time periods may be caused by two factors-different levels of label noise and/or different signal-to-noise ratios. Meanwhile, degradation in performance when transferring a classifier from one time period to another may be caused by three possible factorslabel noise in the source time period, label noise in the target time period, or distributional shift between the time periods. Although isolated differences in predictive performance in a longitudinal setting do not implicate a single causal factor, multiple comparisons taken together may allow us to reason about underlying changes in the data.

This logic guides our search for evidence in support of the hypothesis that mental health annotations cannot be treated as fixed attributes. We consider a standard longitudinal domain transfer setup (Huang and Paul, 2019), chunking the CLPsych dataset into three discrete three-year periods¹ (2012–2015, 2015–2018, 2018–2021) and evaluating within- and between-time-period predictive performance for all available pairs. We use Monte Carlo Cross Validation (Xu and Liang, 2001) to obtain estimates of predictive generalization, chosen over alternative protocols that would be unreliable given the limited sample size of the updated CLPsych dataset (Varoquaux, 2018).

Each iteration of the cross validation procedure (1,000 total) begins by randomly splitting individuals into a 60/40 train/test split, with control and depression groups demographically aligned² using propensity scores (Imbens and Rubin, 2015). To control for differences in data availability between time periods, we not only constrain the sampling process such that splits have an *equal class balance*, but also that individual-level representations are constructed using an *equal document history size* (250 randomly-sampled posts from each time period). A single binary logistic regression classi-

2021
88,.70)
7,.69)
57,.69)

Table 2: Average test-set area under the curve (AUC) and 95% confidence intervals across 1,000 Monte Carlo Cross Validation iterations. Within-time-period performance is significantly higher around the original disclosure window than in subsequent time periods.

fier provided with document-term TF-IDF representations (Baeza-Yates et al., 1999) is fit for each time period using data from individuals in the training set. Each classifier is applied to all three time periods, evaluating performance using individuals in the sampled test set.

Results. We report the average test set area under the curve (AUC) and 95% confidence intervals for each discrete time period pairing in Table 2. Focusing first on within-time-period performance (top left to bottom right diagonal), we find that within-time-period performance is significantly higher in the dataset's original time period (2012-2015) than within subsequent time periods. This holds true even when running experiments only with individuals that have sufficiently-sized post histories in the new time periods, demonstrating that the outcome is not an artifact of survivor bias. At a high level, the differences in within-time-period performance suggest that either label noise has increased or that the signal-to-noise ratio has decreased over time.

Unfortunately, examination of between-timeperiod generalization does not conclusively resolve which of these two factors are responsible for the variation. Focusing first on models trained using data from older time periods (top right triangle), we do not observe any significant difference in predictive performance compared to the benchmarks established by models trained and deployed during the same time period. This serves as a contrast to models deployed on older data (bottom left triangle), where we note that classifiers trained on both of the new time periods incur a loss when being applied to the original CLPsych dataset time period. Interestingly, the absolute differences in performance are minimal. We note that the coefficients of the logistic regression classifiers from each independent time period exhibit significantly positive Pearson correlations, ranging from 0.47 to

¹Time periods were chosen to maximize the number of discrete windows while ensuring enough posts were available to construct informative individual-level representations.

²Aligned on gender and age dimensions.

0.52, and in turn promote stable performance.

Discussion. Although these experiments have not conclusively answered our primary research question regarding longitudinal label validity, they have provided evidence that not all time periods of data are equally informative for training a robust depression classifier. Critically, these results suggest that practitioners cannot assume it better to train a depression classifier using new data, which may be more relevant to their deployment scenario, if it means potentially compromising the temporal relevance of the original ground truth annotations.

What remains to be understood is *why* the predictive task appears to become more difficult in the updated time periods at a statistically significant level, but not one that would necessarily raise immediate concerns to a practitioner. Had underlying dynamics significantly changed since the original data collection period, we would have expected to see a more dramatic loss in predictive performance. Has the mental health status for these individuals genuinely remained static, or is there a spurious confound in the data inflating our performance estimates?

5 Interpreting Model Performance

We attempt to better understand the variation in predictive performance estimated above by comparing language within the updated dataset to the original CLPsych sample. In particular, we adopt a mixed methods approach that allows us to estimate changes in the proportion of depression labels which remain relevant in the updated dataset, and to qualitatively summarize drivers of model decision-making across time periods. We support our analysis by manually coding content-related motifs within a large sample of document histories in the updated dataset, focusing primarily on criteria for diagnosing depression as defined within the DSM-5 (APA, 2013). We draw inspiration from the growing literature on "train-set debugging" (Koh and Liang, 2017; Han et al., 2020), which leverages instance attribution and other diagnostics to succinctly interpret the relationship between training data, learned model parameters, and downstream predictions.

Methods. An annotator is presented with up to 30 anonymized tweets made by a single individual during one of the time periods and asked to indicate whether the individual exhibits evidence of depression. The annotator must mark one of

four options — Uncertain, No Evidence, Some Evidence (Moderate Confidence), Strong Evidence (High Confidence). Explicit disclosures of a depression diagnosis and references to living with depression are automatically assigned to the Strong Evidence category. Otherwise, the annotator is instructed to indicate their confidence based on the nine DSM-5 criteria for diagnosing depression (APA, 2013) and their prior knowledge regarding the presentation of mental health conditions within social media. If at least some evidence of a depression diagnosis is indicated, the annotator is asked to identify whether the depression appears to be in remission (e.g., discussion of overcoming depression). They are also asked to indicate which DSM-5 criteria and/or prior knowledge was used to inform their decision, along with any other notable thematic content.

Our goal of this analysis is *not* to make diagnostic claims regarding the mental health status of individuals in our dataset, but rather to broadly understand what the statistical classifiers are learning. Accordingly, tweets presented to the annotator are those which had the largest positive effect on the classifier's estimated probability of depression, as measured by their influence on user-level predictions within a given time period τ . Formally, we define the influence of a tweet I(x) amongst a set of tweets $x \in X_{\tau}$ as follows:

$$I(x) = \sum_{k=1}^{K} P_{k,\tau}(y = 1|X_{\tau}) - P_{k,\tau}(y = 1|X_{\tau}^{\neg x})$$

where $P_{k,\tau}(\cdot)$ is the probability of depression estimated by a classifier trained on the k-th random sample of data from time period τ , out of K total samples. As was the case in the classification experiments above, each training sample contains 60% of the available data, with the learned classifiers only being applied to the remaining 40% of individuals at each iteration. We refrain from filtering mental health related tweets and those containing explicit diagnosis disclosures, as the goal in this experiment is not to quantify predictive ability, but rather to identify evidence of depression over time. Note that we control for distributional shift over time by estimating influence using a model trained during the time period in which a tweet was posted.

Data. A total of 300 individuals (574 total instances) were selected randomly for annotation. One author, a doctoral student in computer science with multiple years experience working with the

CLPsych dataset, was responsible for all coding. They consulted one additional co-author, an expert in computational modeling of social media and mental health, to develop a common mental model for identifying DSM-5 criteria and other common linguistic motifs in the text. During a pilot round of coding, 16 thematic patterns were identified within the annotated instances to complement the original DSM-5 criteria. Exemplary tweets (paraphrased non-trivially to preserve anonymity (Ayers et al., 2018)) for each of the DSM-5 criteria and alternative thematic categories are provided in Appendix A.3. A breakdown of annotation results is presented in the Table 3. We provide a distribution of the top 20 most common evidence categories amongst individuals who displayed at least some evidence of a depression diagnosis in Appendix

Reliability. Two non-authors with a background in computational psychology independently annotated a subsample of the coded instances to assess the primary coder's reliability. Agreement regarding whether an individual exhibits evidence of depression was fair to moderate; we observe Krippendorff's α measures of 0.438 and 0.499 for the four-class (Uncertain, No Evidence, Some Evidence, Strong Evidence) and three-class (Uncertain, No Evidence, Some or Strong Evidence) scenarios, respectively (Krippendorff, 2011). Agreement regarding remission status varied significantly between pairs of annotators and was generally weaker than agreement regarding evidence of depression ($\alpha = 0.356$). We include an analysis of the disagreements in Appendix A.2 to better contextualize observations from the primary coder's annotations. Succinctly, we identify two reasons for the variation: 1) each annotator's propensity to select the "Uncertain" category, and 2) each annotator's sensitivity to displays of emotion as an indicator of depression.

5.1 What proportion of labels in the updated sample remain relevant?

In line with underlying clinical knowledge regarding the dynamic nature of depression, we observe a significant decrease in linguistic evidence of depression over the course of time. Roughly 76% of individuals in the original depression group displayed at least some clear evidence of a depression diagnosis during the first time period (2012-2015), in comparison to 45% and 39% of individuals in

	Dates	Total	Some Evi.	Strong Evi.	Not Active
•	2012-2015	83	15	3	1
Con	2015-2018	50	10	2	0
0	2018-2021	40	5	0	0
	2012-2015	215	164	136	10
Dep	2015-2018	107	49	28	2
	2018-2021	79	31	16	1

Table 3: Breakdown of coding labels as a function of time period and labels from the original CLPsych dataset. Clinically aligned evidence of a depression diagnosis becomes less prevalent over time.

the 2015-2018 and 2018-2021 time periods, respectively. Across all time periods, only a small number of affirmative instances of depression appear to be in remission. That said, the non-zero level of inactive depression annotations in the original time period highlights an important consideration for practitioners who would like to leverage disclosure-based mechanisms to annotate mental health data moving forward.

The presence of evidence for a depression diagnosis in a subset of the original control group is quite striking. Other studies have raised questions regarding the possible risk of introducing such label noise when curating a control group using a random sampling protocol (Wolohan et al., 2018), though none have provided tangible evidence of this contamination to the best of our knowledge. We see that approximately 4% of individuals in the control group display strong evidence of a depression diagnosis within the original time period. Although relatively small, it is an important reminder of the pitfalls of random control group sampling for health-related social media modeling tasks.

Discussion. The decrease in evidence of a depression diagnosis over time lends support to the introduction of label noise in the updated dataset. Furthermore, it would explain the decrease in predictive performance observed in our previous classification experiments. However, the proportional drop in evidence of a depression diagnosis over time appears too large given the relatively minor reduction in classification accuracy.

We identify two possible explanations for this inconsistency. First, we recognize the possibility that our annotation procedure is insufficient to provide an annotator with appropriate information and comprehensive criteria for indicating evidence of

a depression diagnosis. Only a small subset of an individual's entire post history is displayed to the annotator, a subset chosen using an inherently errorprone statistical ranking method. It is possible that stronger indicators of a depression diagnosis lie outside the 30-tweet sample size window for some individuals. Moreover, the annotator was instructed to rely predominantly on DSM-5 criteria to inform their decision, though several prior computational studies have shown language informative of depression may stray from explicit diagnostic criteria and be difficult for humans to recognize altogether (e.g., increased personal pronoun usage (Holtzman et al., 2017)).

More concerning is the possible presence of non-trivial confounds introduced by the original dataset's sampling/annotation procedure which may artificially inflate predictive performance estimates. Similar types of bias have been identified in prior work when attempting to transfer statistical mental health models trained using proxy-based annotations to new populations of individuals (e.g., demographics, patient populations) (Ernala et al., 2019; Aguirre et al., 2021). Although sampling-based artifacts may be causally-related to the original diagnosis disclosure (e.g., a coping mechanism that becomes a hobby, heightened levels of neuroticism), they may be serve as a red herring in place of primary indicators of depression.

5.2 Do presentations of depression provide evidence of sampling-related confounds?

Personality-related attributes are prominent features in all periods of the updated dataset. For example, indications of a depressed and/or irritable mood were the most common form of evidence in support of an individual having a depression diagnosis. In many cases, anger and irritation were displayed in the form of interpersonal confrontation (passively and actively) with other Twitter profiles. Negative emotions such as loneliness, fear, and existential dread were also displayed readily amongst those showing signs of a depression diagnosis. This result aligns with knowledge regarding the relationship between personality and depression, with elevated levels of neuroticism (negative affectivity and vulnerability to stress) being common in those living with depression (Bagby et al., 2008; Lahey, 2009; Bondy et al., 2021). Although etiologically relevant, this heightened level of emotional affect emerges as one possible artifact which

may confound displays of depression and serve as a nuisance variable in linguistic models of the condition (Tackman et al., 2019).

We also found it common for individuals to mention comorbid psychiatric conditions—such as obsessive compulsive disorder, bipolar disorder, and general anxiety. Many of these conditions share similar underlying symptoms and causes with depressive disorders (Franklin and Zimmerman, 2001; Goodwin, 2015), but tend to assume a different temporal profile (Schoevers et al., 2005). The significant overlap often makes it difficult for trained physicians to properly diagnose individuals (Bowden, 2001) and for language-based algorithms to achieve appropriate discriminative sensitivity (Ive et al., 2018). We recognize the possibility that these comorbid conditions are active during the updated time periods for some individuals and may assume a proxy role in place of depression.

Although not captured by any single evidence category in isolation, there emerged a distinct propensity for "oversharing" amongst individuals from the original dataset's depression group. More specifically, we identified ample discussion of topics that are typically considered socially inappropriate in public discourse spaces (e.g., sexual activity, familial conflict, use of controlled substances). On one hand, this is an interesting finding given that individuals living depression often demonstrate lower levels of emotional self-disclosure (Wei et al., 2005; Kahn and Garrison, 2009). On the other hand, we note that prior work in clinical psychology has recognized a similar propensity for depressed and anxious individuals to engage in oversharing within social media (Radovic et al., 2017; Law et al., 2020).

The theory behind the latter is that social media offers an opportunity to discuss the oft stigmatized challenges of mental health (Betton et al., 2015) and increase feelings of connectedness in a less personal environment (Luo and Hancock, 2020). With this in mind, perhaps it is not surprising that those who have openly disclosed their experience with depression also feel comfortable discussing the aforementioned "taboo" topics. Nonetheless, this personal comfort remains relatively unique amongst the larger social media population. The unfortunate effect of this nuance is that it transforms the primary depression inference task into, essentially, a topic-classification task.

Discussion. Our analysis affirms what other recent studies on proxy-based mental health annota-

tions have claimed — individuals who disclose a mental health condition systematically differ from the larger population of individuals living with that condition (Ernala et al., 2019; Saha et al., 2021). As a research community, we must be careful to disambiguate 1) training a language classifier to identify individuals who live with a mental health condition, and 2) training a language classifier to identify individuals who live with a mental health condition and disclose their diagnosis. Inappropriately equating the two creates an opportunity to erroneously estimate population-level dynamics (Amir et al., 2019) and ignore underrepresented voices from communities who tend to possess conservative ideologies regarding mental health (Loveys et al., 2018; Aguirre et al., 2021).

6 Discussion

Demand for computational methods to quantify mental health dynamics within social media data is at an all time high (Galea et al., 2020). However, the potential impact of these methods remains bounded by the robustness of datasets used for their development. Spanning nearly a decade of online activity, our study uniquely identifies evidence of these limitations as they currently manifest in non-clinically derived mental health social media datasets. This evidence leads us to offer three recommendations for enhancing data curation and model evaluation.

Annotate Diagnosis Date & Comorbidities. We identified several instances within our dataset where a diagnosis disclosure was made in reference to a condition that had since entered remission. In other cases, depression diagnoses were either supplanted by or augmented with alternative psychiatric diagnoses. Indicators regarding the time a diagnosis was made, many of which can be identified using inexpensive algorithms (MacAvaney et al., 2018), can provide important signal regarding the temporal relevance of a psychiatric diagnosis. Meanwhile, inclusion of comorbidities may provide researchers an opportunity to model psychiatric heterogeneity (Arseniev-Koehler et al., 2018) and interpret longitudinal generalization.

Sample Control Groups using Propensity Matching. Control group selection is influential in both training and evaluation of statistical models of mental health (Pirina and Çöltekin, 2018). Prior work has leveraged a myriad of criteria to match individuals who have disclosed a psychiatric di-

agnosis with suitable counterparts—demographics (Coppersmith et al., 2014), online behavior (Cohan et al., 2018), and language (De Choudhury et al., 2016). Though use of inconsistent matching criteria is less than ideal, the absence of any protocol is potentially more problematic (Shen et al., 2018; Wolohan et al., 2018). We recommend practitioners leverage propensity-based matching (Imbens and Rubin, 2015) to reduce the effect of self-disclosure biases (e.g., personality, interests, demographics). In addition to the aforementioned dimensions, researchers may augment their criteria using classifiers to infer relevant latent attributes (Preoţiuc-Pietro et al., 2015) or neural models to derive user-level embeddings (Amir et al., 2017).

Identify and Filter Sampling Biases. Our analysis benefited from context that emerged when attempting to train classifiers that generalize over long time periods. However, access to supplementary data is not necessary to understand whether artifacts may exist in a dataset. Algorithmic approaches, such as those from Le Bras et al. (2020), may be used to identify instances containing spurious correlations. These approaches should be used to augment insights derived from manual annotation and review. We found our technique for ranking the influence of individual posts on user-level predictions began yielding insights after only a few dozen examples, though alternative ranking methodologies are available (Uban et al., 2021). Outcomes should be used to inform preprocessing decisions, construct fair evaluations (Poliak et al., 2018), and inform the description of a dataset within documentation/datasheets (Gebru et al., 2021).

6.1 Limitations and Qualifiers

Though our analysis identified data attributes that may inhibit statistical generalization, we also found evidence in support of the validity of self-disclosed diagnoses for annotating mental health status. The majority of individuals within the CLPsych dataset's original time window showed clear evidence of depression that aligns with clinical criteria. Many of these indicators remained stable over the course of time. Moreoever, the 2015 CLPsych Shared Task dataset is just one of many resources in this research community, all of which are likely to exhibit varying degrees of noise depending on their respective sampling protocols. Conclusive statements regarding the validity of self-disclosed

diagnoses require evidence from multiple social media platforms, cultural groups, and time periods.

7 Ethical Considerations

Ethical challenges emerging from use of public social media data to analyze an individual's mental health have been examined extensively by members of both computational and clinical/public health communities (Conway and O'Connor, 2016; Chancellor et al., 2019). Privacy-related concerns are the most poignant for our study, which relies both on de-anonymizing records from a vulnerable population and manually reviewing/analyzing individual posts.

Indeed, many individuals who publicly discuss their mental health or disclose a psychiatric condition within social media admit that they worry about harmful repercussions of sharing such sensitive information with the public (Ford et al., 2019; Naslund and Aschbrenner, 2019). Primary fears include risking occupational stability, damaging interpersonal relationships, and being subjected to hostile communications. Whether potential positive outcomes (e.g., development of systems for recommending mental health care, fiduciary aid to address population-level crises) offset these threats remains largely dependent on an individual's personal life experience. For example, psychiatric patients have expressed stronger approval toward analysis of their social media than members of the general public (Mikal et al., 2017). The same holds true amongst younger individuals (Naslund and Aschbrenner, 2019).

Recognizing these viewpoints, we are careful to mitigate privacy-related risks to the greatest extent possible given our primary research aim. For example, account identifiers distributed within the 2015 CLPsych Shared Task dataset are de-anonymized only temporarily to link updated records with existing post histories. We also redact account handles and URLs from the text analyzed during our manual coding procedure (§5). In line with protocols enumerated by Benton et al. (2017a), all data is stored on a remote server and secured using OSlevel group permissions. We perform our analysis under the external guidance of clinical psychologists and psychiatrists. Our study is also reviewed by our Institutional Review Board (IRB), obtaining exempt status under 45 CFR §46.104.

Critically, our intention is not to develop a publicfacing system for algorithmic analysis of mental health. Rather, our goal is to evaluate the validity of an existing and widely-adopted data curation practice (Chancellor and De Choudhury, 2020; Harrigian et al., 2021). Failure to comprehensively understand biases that arise under this methodology can have severe detrimental effects in downstream systems. In the case of estimating populationlevel health trends, for instance, we have already seen machine learning classifiers produce outcomes that are inconsistent across computational studies (Wolohan, 2020; Biester et al., 2021; Harrigian and Dredze, 2022) and in conflict with traditional measurement techniques (Amir et al., 2019). Continuing to pursue this line of research without questioning the validity of its underlying data has the potential to irreparably damage the public's trust in this domain, and worse, enable ill-informed decision making in highly-sensitive circumstances.

Acknowledgements

We thank Ayah Zirikly and Carlos Aguirre for contributing annotations to use for evaluating interrater reliability. We also thank the anonymous reviewers for providing additional clinical grounding of our study and highlighting opportunities to improve our technical approach.

References

Oshin Agarwal and Ani Nenkova. 2021. Temporal effects on pre-trained models for language processing tasks. *arXiv preprint arXiv:2111.12790*.

Carlos Aguirre, Keith Harrigian, and Mark Dredze. 2021. Gender and racial fairness in depression research using social media. In 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume.

Silvio Amir, Glen Coppersmith, Paula Carvalho, Mário J Silva, and Bryon C Wallace. 2017. Quantifying mental health from social media with neural user embeddings. In *Machine Learning for Health-care Conference*. PMLR.

Silvio Amir, Mark Dredze, and John W Ayers. 2019. Mental health surveillance over social media with digital cohorts. In Sixth Workshop on Computational Linguistics and Clinical Psychology.

Jules Angst, Alex Gamma, Wulf Rössler, Vladeta Ajdacic, and Daniel N Klein. 2009. Long-term depression versus episodic major depression: results from the prospective zurich study of a community sample. *Journal of affective disorders*, 115(1-2).

APA. 2013. Diagnostic and statistical manual of mental disorders: Dsm-5. *Arlington, VA*.

- Alina Arseniev-Koehler, Sharon Mozgai, and Stefan Scherer. 2018. What type of happiness are you looking for?-a closer look at detecting mental health from language. In *Fifth Workshop on Computational Linguistics and Clinical Psychology: From Keyboard to Clinic*.
- John W Ayers, Theodore L Caputi, Camille Nebeker, and Mark Dredze. 2018. Don't quote me: reverse identification of research participants in social media studies. NPJ digital medicine.
- Ricardo Baeza-Yates, Berthier Ribeiro-Neto, et al. 1999. Modern information retrieval. ACM press New York.
- R Michael Bagby, Lena C Quilty, and Andrew C Ryder. 2008. Personality and depression. *The Canadian Journal of Psychiatry*.
- Adrian Benton, Glen Coppersmith, and Mark Dredze. 2017a. Ethical research protocols for social media health research. In *First ACL Workshop on Ethics in Natural Language Processing*.
- Adrian Benton, Margaret Mitchell, and Dirk Hovy. 2017b. Multitask learning for mental health conditions with limited social media data. In 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers.
- Victoria Betton, Rohan Borschmann, Mary Docherty, Stephen Coleman, Mark Brown, and Claire Henderson. 2015. The role of social media in reducing stigma and discrimination. *The British Journal of Psychiatry*.
- Laura Biester, Katie Matton, Janarthanan Rajendran, Emily Mower Provost, and Rada Mihalcea. 2021. Understanding the impact of covid-19 on online mental health forums. ACM Transactions on Management Information Systems (TMIS).
- Erin Bondy, David AA Baranger, Jared Balbona, Kendall Sputo, Sarah E Paul, Thomas F Oltmanns, and Ryan Bogdan. 2021. Neuroticism and reward-related ventral striatum activity: Probing vulnerability to stress-related depression. *Journal of Abnormal Psychology*, 130(3):223.
- Charles L Bowden. 2001. Strategies to reduce misdiagnosis of bipolar depression. Psychiatric Services.
- Luca Canzian and Mirco Musolesi. 2015. Trajectories of depression: unobtrusive monitoring of depressive states by means of smartphone mobility traces analysis. In 2015 ACM international joint conference on pervasive and ubiquitous computing.
- Stevie Chancellor, Michael L Birnbaum, Eric D Caine, Vincent MB Silenzio, and Munmun De Choudhury. 2019. A taxonomy of ethical tensions in inferring mental health states from social media. In conference on fairness, accountability, and transparency.

- Stevie Chancellor and Munmun De Choudhury. 2020. Methods in predictive techniques for mental health status on social media: a critical review. *NPJ digital medicine*, 3(1).
- Stevie Chancellor, Zhiyuan Lin, Erica L Goodman, Stephanie Zerwas, and Munmun De Choudhury. 2016. Quantifying and predicting mental illness severity in online pro-eating disorder communities. In 19th ACM Conference on Computer-Supported Cooperative Work & Social Computing.
- Daejin Choi, Steven A Sumner, Kristin M Holland, John Draper, Sean Murphy, Daniel A Bowen, Marissa Zwald, Jing Wang, Royal Law, Jordan Taylor, et al. 2020. Development of a machine learning model using multiple, heterogeneous data sources to estimate weekly us suicide fatalities. *JAMA network open*.
- Arman Cohan, Bart Desmet, Andrew Yates, Luca Soldaini, Sean MacAvaney, and Nazli Goharian. 2018. Smhd: a large-scale resource for exploring online language usage for multiple mental health conditions. In 27th International Conference on Computational Linguistics. ACL.
- Jacob Cohen. 1960. A coefficient of agreement for nominal scales. Educational and psychological measurement, 20(1):37–46.
- Stephan Collishaw, Barbara Maughan, Robert Goodman, and Andrew Pickles. 2004. Time trends in adolescent mental health. *Journal of Child Psychology and psychiatry*, 45(8).
- Mike Conway and Daniel O'Connor. 2016. Social media, big data, and mental health: current advances and ethical implications. *Current opinion in psychology*.
- Glen Coppersmith, Mark Dredze, and Craig Harman. 2014. Quantifying mental health signals in twitter. In workshop on computational linguistics and clinical psychology: From linguistic signal to clinical reality.
- Glen Coppersmith, Mark Dredze, Craig Harman, Kristy Hollingshead, and Margaret Mitchell. 2015. Clpsych 2015 shared task: Depression and ptsd on twitter. In 2nd Workshop on Computational Linguistics and Clinical Psychology: From Linguistic Signal to Clinical Reality.
- Cassandra R Davis, Jevay Grooms, Alberto Ortega, Joaquin Alfredo-Angel Rubalcaba, and Edward Vargas. 2021. Distance learning and parental mental health during covid-19. Educational Researcher.
- Munmun De Choudhury and Sushovan De. 2014. Mental health discourse on reddit: Self-disclosure, social support, and anonymity. In *Eighth international AAAI conference on weblogs and social media*.

- Munmun De Choudhury, Emre Kiciman, Mark Dredze, Glen Coppersmith, and Mrinal Kumar. 2016. Discovering shifts to suicidal ideation from mental health content in social media. In 2016 CHI conference on human factors in computing systems.
- Orianna DeMasi, Konrad Kording, and Benjamin Recht. 2017. Meaningless comparisons lead to false optimism in medical machine learning. *PloS one*.
- Sindhu Kiranmai Ernala, Michael L Birnbaum, Kristin A Candan, Asra F Rizvi, William A Sterling, John M Kane, and Munmun De Choudhury. 2019. Methodological gaps in predicting mental health states from social media: triangulating diagnostic signals. In 2019 chi conference on human factors in computing systems.
- Alex Fine, Patrick Crutchley, Jenny Blase, Joshua Carroll, and Glen Coppersmith. 2020. Assessing population-level symptoms of anxiety, depression, and suicide risk in real time using nlp applied to social media data. In Fourth Workshop on Natural Language Processing and Computational Social Science.
- Elizabeth Ford, Keegan Curlewis, Akkapon Wongkoblap, and Vasa Curcin. 2019. Public opinions on using social media content to identify users with depression and target mental health care advertising: mixed methods survey. *JMIR mental health*.
- C Laurel Franklin and Mark Zimmerman. 2001. Posttraumatic stress disorder and major depressive disorder: Investigating the role of overlapping symptoms in diagnostic comorbidity. *The Journal of nervous* and mental disease.
- Benoît Frénay and Michel Verleysen. 2013. Classification in the presence of label noise: a survey. *IEEE transactions on neural networks and learning systems*.
- Sandro Galea, Raina M Merchant, and Nicole Lurie. 2020. The mental health consequences of covid-19 and physical distancing: the need for prevention and early intervention. *JAMA internal medicine*.
- Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach, Hal Daumé Iii, and Kate Crawford. 2021. Datasheets for datasets. *Communications of the ACM*.
- Guy M Goodwin. 2015. The overlap between anxiety, depression, and obsessive-compulsive disorder. *Dialogues in clinical neuroscience*.
- Sharath Chandra Guntuku, David B Yaden, Margaret L Kern, Lyle H Ungar, and Johannes C Eichstaedt. 2017. Detecting depression and mental illness on social media: an integrative review. Current Opinion in Behavioral Sciences.

- Xiaochuang Han, Byron C Wallace, and Yulia Tsvetkov. 2020. Explaining black box predictions and unveiling data artifacts through influence functions. arXiv preprint arXiv:2005.06676.
- Keith Harrigian, Carlos Aguirre, and Mark Dredze. 2020. Do models of mental health based on social media data generalize? In 2020 conference on empirical methods in natural language processing: findings.
- Keith Harrigian, Carlos Aguirre, and Mark Dredze. 2021. On the state of social media data for mental health research. In Seventh Workshop on Computational Linguistics and Clinical Psychology: Improving Access.
- Keith Harrigian and Mark Dredze. 2022. The problem of semantic shift in longitudinal monitoring of social media. *Proceedings of the 14th ACM Web Science Conference*.
- Nicholas S Holtzman et al. 2017. A meta-analysis of correlations between depression and first person singular pronoun use. *Journal of Research in Personality*.
- Xiaolei Huang and Michael Paul. 2019. Neural temporality adaptation for document classification: Diachronic word embeddings and domain adaptation models. In 57th Annual Meeting of the Association for Computational Linguistics.
- Guido W Imbens and Donald B Rubin. 2015. Causal inference in statistics, social, and biomedical sciences. Cambridge University Press.
- Julia Ive, George Gkotsis, Rina Dutta, Robert Stewart, and Sumithra Velupillai. 2018. Hierarchical neural model with attention mechanisms for the classification of social media text related to mental health. In Fifth Workshop on Computational Linguistics and Clinical Psychology: From Keyboard to Clinic. Association for Computational Linguistics.
- Judith Johnson, Louise H Hall, Kathryn Berzins, John Baker, Kathryn Melling, and Carl Thompson. 2018. Mental healthcare staff well-being and burnout: A narrative review of trends, causes, implications, and recommendations for future interventions. *International journal of mental health nursing*.
- Sheri L Johnson and Andrzej Nowak. 2002. Dynamical patterns in bipolar depression. *Personality and Social Psychology Review*.
- Jeffrey H Kahn and Angela M Garrison. 2009. Emotional self-disclosure and emotional avoidance: Relations with symptoms of depression and anxiety. *Journal of counseling psychology*, 56(4):573.
- Pang Wei Koh and Percy Liang. 2017. Understanding black-box predictions via influence functions. In *International conference on machine learning*, pages 1885–1894. PMLR.

- Klaus Krippendorff. 2004. Reliability in content analysis: Some common misconceptions and recommendations. *Human communication research*, 30(3):411–433.
- Klaus Krippendorff. 2011. Computing krippendorff's alpha-reliability.
- Benjamin B Lahey. 2009. Public health significance of neuroticism. *American Psychologist*, 64(4):241.
- Danielle M Law, Jennifer D Shapka, and Rebecca J Collie. 2020. Who might flourish and who might languish? adolescent social and mental health profiles and their online experiences and behaviors. *Human Behavior and Emerging Technologies*, 2(1):82–92.
- Ronan Le Bras, Swabha Swayamdipta, Chandra Bhagavatula, Rowan Zellers, Matthew Peters, Ashish Sabharwal, and Yejin Choi. 2020. Adversarial filters of dataset biases. In *International Conference on Machine Learning*. PMLR.
- Andrew Lee, Jonathan K Kummerfeld, Larry An, and Rada Mihalcea. 2021. Micromodels for efficient, explainable, and reusable systems: A case study on mental health. In *Findings of the Association for Computational Linguistics: EMNLP 2021*.
- David E Losada, Fabio Crestani, and Javier Parapar. 2017. erisk 2017: Clef lab on early risk prediction on the internet: experimental foundations. In *International Conference of the Cross-Language Evaluation Forum for European Languages*. Springer.
- Kate Loveys, Jonathan Torrez, Alex Fine, Glen Moriarty, and Glen Coppersmith. 2018. Cross-cultural differences in language markers of depression online. In *Proceedings of the fifth workshop on computational linguistics and clinical psychology: from keyboard to clinic*, pages 78–87.
- Marco Lui and Timothy Baldwin. 2012. langid. py: An off-the-shelf language identification tool. In *ACL 2012 system demonstrations*.
- Mufan Luo and Jeffrey T Hancock. 2020. Self-disclosure and social media: motivations, mechanisms and psychological well-being. *Current Opinion in Psychology*, 31:110–115.
- Sean MacAvaney, Bart Desmet, Arman Cohan, Luca Soldaini, Andrew Yates, Ayah Zirikly, and Nazli Goharian. 2018. Rsdd-time: Temporal annotation of self-reported mental health diagnoses. In Fifth Workshop on Computational Linguistics and Clinical Psychology: From Keyboard to Clinic.
- Sean Macavaney, Anjali Mittu, Glen Coppersmith, Jeff Leintz, and Philip Resnik. 2021. Community-level research on suicidality prediction in a secure environment: Overview of the clpsych 2021 shared task. In Seventh Workshop on Computational Linguistics and Clinical Psychology: Improving Access.

- Jude Mikal, Samantha Hurst, and Mike Conway. 2017. Investigating patient attitudes towards the use of social media data to augment depression diagnosis and treatment: a qualitative study. In *fourth workshop on computational linguistics and clinical psychology—from linguistic signal to clinical reality*.
- Jeanne Miranda, Lonnie R Snowden, and Rupinder K Legha. 2020. Policy effects on mental health status and mental health care disparities. In *The palgrave handbook of American mental health policy*. Springer.
- Danielle Mowery, Hilary Smith, Tyler Cheney, Greg Stoddard, Glen Coppersmith, Craig Bryan, and Mike Conway. 2017. Understanding depressive symptoms and psychosocial stressors on twitter: a corpus-based study. *Journal of medical Internet research*, 19(2).
- Aaron Mueller, Zach Wood-Doughty, Silvio Amir, Mark Dredze, and Alicia Lynn Nobles. 2021. Demographic representation and collective storytelling in the me too twitter hashtag activism movement. ACM on Human-Computer Interaction.
- John A Naslund and Kelly A Aschbrenner. 2019. Risks to privacy with use of social media: understanding the views of social media users with serious mental illness. *Psychiatric services*.
- Alicia L Nobles, Jeffrey J Glenn, Kamran Kowsari, Bethany A Teachman, and Laura E Barnes. 2018. Identification of imminent suicide risk among young adults using text messages. In 2018 CHI Conference on Human Factors in Computing Systems.
- Inna Pirina and Çağrı Çöltekin. 2018. Identifying depression on reddit: The effect of training data. In 2018 EMNLP Workshop SMM4H: The 3rd Social Media Mining for Health Applications Workshop & Shared Task.
- Adam Poliak, Jason Naradoesky, Aparajita Haldar, Rachel Rudinger, and Benjamin Van Durme. 2018. Hypothesis only baselines in natural language inference. In Seventh Joint Conference on Lexical and Computational Semantics.
- Daniel Preoţiuc-Pietro, Johannes Eichstaedt, Gregory Park, Maarten Sap, Laura Smith, Victoria Tobolsky, H Andrew Schwartz, and Lyle Ungar. 2015. The role of personality, age, and gender in tweeting about mental illness. In 2nd workshop on computational linguistics and clinical psychology: From linguistic signal to clinical reality.
- Ana Radovic, Theresa Gmelin, Bradley D Stein, and Elizabeth Miller. 2017. Depressed adolescents' positive and negative use of social media. *Journal of adolescence*, 55:5–15.
- Guozheng Rao, Yue Zhang, Li Zhang, Qing Cong, and Zhiyong Feng. 2020. Mgl-cnn: A hierarchical posts representations model for identifying depressed individuals in online forums. *IEEE Access*.

- Andrew G Reece, Andrew J Reagan, Katharina LM Lix, Peter Sheridan Dodds, Christopher M Danforth, and Ellen J Langer. 2017. Forecasting the onset and course of mental illness with twitter data. *Scientific reports*.
- Farig Sadeque, Dongfang Xu, and Steven Bethard. 2018. Measuring the latency of depression detection in social media. In *Eleventh ACM International Conference on Web Search and Data Mining*.
- Koustuv Saha, John Torous, Emre Kiciman, and Munmun De Choudhury. 2021. Understanding side effects of antidepressants: Large-scale longitudinal study on social media data. *JMIR mental health*.
- Maarten Sap, Gregory Park, Johannes Eichstaedt, Margaret Kern, David Stillwell, Michal Kosinski, Lyle Ungar, and H Andrew Schwartz. 2014. Developing age and gender predictive lexica over social media. In 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP).
- Robert A Schoevers, DJH Deeg, W Van Tilburg, and ATF Beekman. 2005. Depression and generalized anxiety disorder: co-occurrence and longitudinal patterns in elderly patients. *The American Journal of Geriatric Psychiatry*.
- Tiancheng Shen, Jia Jia, Guangyao Shen, Fuli Feng, Xiangnan He, Huanbo Luan, Jie Tang, Thanassis Tiropanis, Tat Seng Chua, and Wendy Hall. 2018. Cross-domain depression detection via harvesting social media. In 2018 International Joint Conferences on Artificial Intelligence (IJCAI). International Joint Conferences on Artificial Intelligence.
- Masashi Sugiyama and Motoaki Kawanabe. 2012. *Machine learning in non-stationary environments: Introduction to covariate shift adaptation*. MIT press.
- Allison M Tackman, David A Sbarra, Angela L Carey, M Brent Donnellan, Andrea B Horn, Nicholas S Holtzman, To'Meisha S Edwards, James W Pennebaker, and Matthias R Mehl. 2019. Depression, negative emotionality, and self-referential language: A multi-lab, multi-measure, and multi-language-task research synthesis. *Journal of personality and social psychology*, 116(5):817.
- Adam Tsakalidis, Maria Liakata, Theo Damoulas, and Alexandra I Cristea. 2018. Can we assess mental health through social media and smart devices? addressing bias in methodology and evaluation. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*. Springer.
- Ana Sabina Uban, Berta Chulvi, and Paolo Rosso. 2021. Understanding patterns of anorexia manifestations in social media data with deep learning. In Seventh Workshop on Computational Linguistics and Clinical Psychology: Improving Access.
- Gaël Varoquaux. 2018. Cross-validation failure: small sample sizes lead to large error bars. *Neuroimage*, 180.

- Matej Gjurkovic Mladen Karan Iva Vukojevic and Mihaela Bošnjak Jan Šnajder. 2021. Pandora talks: Personality and demographics on reddit. *SocialNLP* 2021.
- Meifen Wei, Daniel W Russell, and Robyn A Zakalik. 2005. Adult attachment, social self-efficacy, self-disclosure, loneliness, and subsequent depression for freshman college students: A longitudinal study. *Journal of counseling psychology*, 52(4):602.
- JT Wolohan. 2020. Estimating the effect of covid-19 on mental health: Linguistic indicators of depression during a global pandemic. In 1st Workshop on NLP for COVID-19 at ACL 2020.
- JT Wolohan, Misato Hiraga, Atreyee Mukherjee, Zeeshan Ali Sayyed, and Matthew Millard. 2018. Detecting linguistic traces of depression in topicrestricted text: Attending to self-stigmatized depression with nlp. In First International Workshop on Language Cognition and Computational Models.
- Akkapon Wongkoblap, Miguel A Vadillo, and Vasa Curcin. 2019. Predicting social network users with depression from simulated temporal data. In *IEEE EUROCON 2019-18th International Conference on Smart Technologies*. IEEE.
- Qing-Song Xu and Yi-Zeng Liang. 2001. Monte carlo cross validation. *Chemometrics and Intelligent Laboratory Systems*.

A Interpreting Model Performance

A.1 Data

Three individuals (one author A_1 , two non-authors B_1, B_2) independently generated the annotations used to facilitate the analysis presented in §5. Statistics presented in the analysis are computed using the author's annotations, while reliability measures are computed using additional annotations from the non-authors. All annotators have several years of experience modeling language within social media to assess mental health, but do not claim to be experts in clinical psychology. Additionally, all annotators have prior experience with the CLPsych 2015 Shared Task data (Coppersmith et al., 2015) — e.g., A_1 and B_1 have worked with the original CLPsych dataset extensively over the prior three years. We include the distribution of instances reviewed by each of our annotators in Table 4.

A.2 Inter-rater Reliability

As a first look into inter-rater reliability, we consider three dimensions of agreement — evidence

Time Period

	2012-2015	2015-2018	2018-2021	Total
$\overline{A_1}$	298	157	119	574
B_1	103	62	40	205
B_2	26	15	12	53

Table 4: Distribution of instances coded by each annotator across the three time periods. Note that the set of instances annotated follows the relationship: $B_2 \subseteq B_1 \subseteq A_1$.

of depression (four-class and three-class)³ and remission status (four-class). We present pairwise annotator agreement matrices for each of these dimensions in Figure 1. We use Cohen's kappa κ to evaluate pairwise annotator agreement (Cohen, 1960) and Krippendorff's alpha α to evaluate multiannotator agreement (Krippendorff, 2011).

We observe fair to moderate agreement for the evidence-of-depression task: $\alpha=0.4376$ and $\alpha=0.4988$ for the four-class and three-class versions, respectively. Meanwhile, agreement on remission status is poor, reflected by a Krippendorff's α of 0.3561. In isolation, these agreement measures would suggest the results of our analysis should be accepted tentatively at best (Krippendorff, 2004). However, we argue these statistics are perhaps a bit conservative and skewed by the small sample size of annotations generated by B_2 . A review of the underlying distributions provides us an opportunity to understand axes of disagreement and, in turn, contextualize the results presented in §5.

As shown in Figure 1, annotator B_2 exhibits a higher propensity to use the "Uncertain" label in the evidence-of-depression tasks compared to annotators A_1 and B_1 . At the same time, while annotator B_2 is more inclined to indicate they are uncertain about an example than annotator A_1 , we note that annotator B_1 appears to have a higher baseline threshold of what constitutes evidence of depression than annotator A_1 . The latter is demonstrated by the fact that nearly all examples marked in the affirmative by B_1 were also marked as such by A_1 , but a large number of examples marked in the affirmative by A_1 were marked as not containing evidence of depression by B_1 .

With respect to the remission status task (bottom subplot of Figure 1), we note that annotator

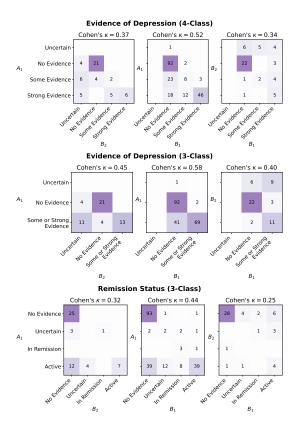


Figure 1: Pairwise agreement matrices for the annotation tasks. Underlying relationships reveal cognitive biases from annotator A_1 that may affect the outcomes presented in §5.

 B_1 is more likely to mark an example as uncertain and more likely to mark an example as being inremission than annotators A_1 and B_2 . Broadly, this distribution highlights the difficulty of distinguishing active cases of clinical depression from prior experiences and lingering effects. It also serves as support for our recommendation in §6 that researchers should attempt to include the time a diagnosis was received by an individual when curating new datasets.

We acquire additional context for our results by examining the distribution of annotations as a function of the original CLPsych labels. Examining the results visualized in Figure 2, we first note that annotator A_1 classifies instances most accurately (under the assumption that ground truth is fixed over time). We believe this outcome to be a result of exposure bias; the annotation task was conducted *after* the completion of several modeling experiments, through which annotator A_1 was uniquely provided an opportunity to learn more about the presentation of depression by individuals

³Note that the three-class evidence-of-depression grouping simply merges the Some Evidence and Strong Evidence categories of the four-class version.

in the 2015 CLPsych Shared Task dataset. We also note the distribution of "Uncertain" decisions from annotator B_2 concentrating within the original depression group. This seems to suggest annotator B_2 adopted a conservative coding approach when presented with instances that contained smaller degrees of evidence, whereas annotators A_1 and B_1 required a lower threshold of evidence to make a decision.

To conclude our reliability analysis, we examine agreement regarding the manner in which each annotator made their decision (i.e., evidence identification). We find that annotators A_1 and B_1 generally identify diagnosis disclosures within the same instances. Annotator B_2 often abstained from making a decision when presented with a disclosure due to uncertainty regarding the subject of the diagnosis. Annotator A_1 also indicated the presence of a depressed and/or irritable mood at a significantly higher rate than the other annotators, seemingly more sensitive to extreme negative emotions than the other annotators.

Discussion. Considering the difficulty of the annotation task, it is perhaps not surprising to have observed less than perfect annotator agreement. Machine learning classifiers often require hundreds of posts to make an accurate estimate of an individual's mental health status, while our annotators were only provided at maximum of 30 posts and encouraged to rely on varying levels of prior knowledge regarding the presentation of depression in social media. Critically, we emphasize that the goal of the analysis presented in §5 is *not* to curate ground truth labels of mental health status or act as clinical experts, but rather to understand biases that may exist in a depression dataset generated using self-disclosed diagnoses. The analysis of inter-rater reliability presented above provides an opportunity to further ground the results discussed in §5 and highlight areas that may benefit from future research.

A.3 Evidence Distribution

We include a breakdown of evidence annotations for individuals displaying some evidence of depression (§5) in Figure 3. Exemplary tweets for each of the evidence categories (paraphrased to maintain anonymity) are provided in Table 5. Both can be found on the following pages.

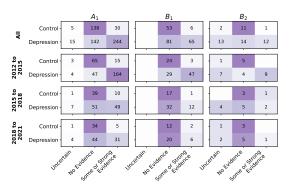


Figure 2: Distribution of annotations for the evidence of depression task (three-class) as a function of the original CLPsych labels. Affirmative evidence becomes less prevalent in the new time periods compared to the original time period for each annotator.

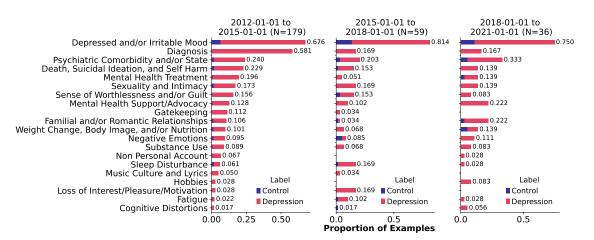


Figure 3: Distribution of evidence amongst individuals indicated as displaying at least some evidence of a depression diagnosis. Depressed and/or irritable mood is consistently the most common type of evidence within each of the three time periods.

Evidence	Exemplary Tweets
Diagnosis Disclosure	"Bipolar disorder and depression. My doctor finally agrees." "I have suffered from depression for several years now"
Depressed & Irritable Mood	"No one ever asks if I'm doing fine."
Loss of Interest/Pleasure/Motivation	"You don't understand what I'm dealing with. Get fucked." "realizing you don't care about the things you used to enjoy" "cant get out of bed today"
Weight, Body Image, & Nutrition	"Not that anyone cares, but I'm almost at my goal weight." "I bought the dress I've always wanted, but still don't feel pretty."
Sleep Disturbance	"I CANT SLEEP. PAIN. JUST LIKE ALWAYS." "Shit! Surviving on only a couple of hours of sleep again :/"
Fatigue	"mentally drained from this pandemic" "This should be effortless but I can't work any harder"
Sense of Worthlessness & Guilt	"when you let someone do anything to you" "It truly is always my fault. I probably suck."
Impaired Thought	"I'm failing my classes because I'm depressed." "at work. cant focus doe"
Death & Self Harm	"My scars are fadedunless you care to look close" "I wish you all never see a loved one fade away."
Cognitive Distortions	"Going to fail this exam. SCREWED." "I always think my bf is going to leave me"
Treatment	"Scared to tell a women that I'm in therapy" "Slowly weaning of the prozac."
Gatekeeping	"depression isn't just a bad day. fuck you all." "LET ME SHOW YOU WANT DEPRESSION IS"
Sexuality and Intimacy	"Who wants to come take some pics of me for only fans?;)" "Every girl should watch porn with their bf"
Negative Emotions	"hi sunshine! Too bad no one to spend today with." "I feel like no one cares even though I know they do"
Coping Strategies	"Have you talked to anyone about it yet?" "Art is always the easiest way to distract me from my anxiety"
Psychiatric Comorbidity & State	"Really stressing today. Lots of built up anger" "I am anorexic and cut myself."
Non-psychiatric Comorbidity	"Could use a little bit of aid #DisabilityAid" "Lots of back pain ruining what should be a beautiful day."
Substance Use	"I really shouldn't be drunk this early." "Weed makes the dreams go away and thats a good thing."
Support & Advocacy	"If I can manage a smile, I believe you can too one day!" "RIP Chester. If you're going through pain, reach out to me."
Personality and Identity	"Girls say they love a man in uniform until they do their job" "Lol grandma still think I'm bringing a boy home"
Music Culture & Lyrics	"#FallingInReverse :D" "Scene doesn't mean emo idiots. I dont want to kill myself."
Familial/Romantic Relationships	"when bae dont answer the phone xx"
Political & Moral Beliefs	"Mom: You'll never lose weight. Me: Is that why dad left?" "look in the mirror if you're not upset a cop can murder" "Trusper will kill up all"
Hobbies	"Trump will kill us all" "Missin the old days when eveyone played Pokemon yellow"
Non-personal Accounts	"Boys that watch the Kardashians. Love." "My life was about to fall apart until I found the Calm app" "Breaking News: 5-alarm fire just outside Tulsa"

Table 5: Exemplary tweets and phrases (modified to preserve anonymity) for each of the 25 evidence categories.

A Interpreting Model Performance

A.1 Data

Three individuals (one author A_1 , two non-authors B_1, B_2) independently generated the annotations used to facilitate the analysis presented in §5. Statistics presented in the analysis are computed using the author's annotations, while reliability measures are computed using additional annotations from the non-authors. All annotators have several years of experience modeling language within social media to assess mental health, but do not claim to be experts in clinical psychology. Additionally, all annotators have prior experience with the CLPsych 2015 Shared Task data (?) — e.g., A_1 and B_1 have worked with the original CLPsych dataset extensively over the prior three years. We include the distribution of instances reviewed by each of our annotators in Table 1.

Time Period

	2012-2015	2015-2018	2018-2021	Total
A_1	298	157	119	574
B_1	103	62	40	205
B_2	26	15	12	53

Table 1: Distribution of instances coded by each annotator across the three time periods. Note that the set of instances annotated follows the relationship: $B_2 \subseteq B_1 \subseteq A_1$.

A.2 Inter-rater Reliability

As a first look into inter-rater reliability, we consider three dimensions of agreement — evidence of depression (four-class and three-class) and remission status (four-class). We present pairwise annotator agreement matrices for each of these dimensions in Figure 1. We use Cohen's kappa κ to evaluate pairwise annotator agreement (?) and Krippendorff's alpha α to evaluate multi-annotator agreement (?).

We observe fair to moderate agreement for the evidence-of-depression task: $\alpha=0.4376$ and $\alpha=0.4988$ for the four-class and three-class versions, respectively. Meanwhile, agreement on remission status is poor, reflected by a Krippendorff's α of 0.3561. In isolation, these agreement measures would suggest the results of our analysis

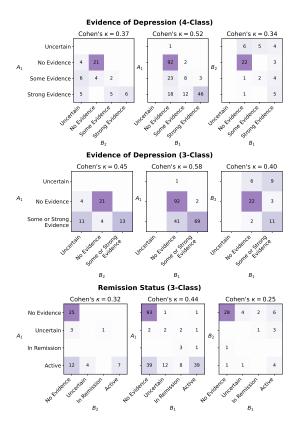


Figure 1: Pairwise agreement matrices for the annotation tasks. Underlying relationships reveal cognitive biases from annotator A_1 that may affect the outcomes presented in §5.

should be accepted tentatively at best (?). However, we argue these statistics are perhaps a bit conservative and skewed by the small sample size of annotations generated by B_2 . A review of the underlying distributions provides us an opportunity to understand axes of disagreement and, in turn, contextualize the results presented in §5.

As shown in Figure 1, annotator B_2 exhibits a higher propensity to use the "Uncertain" label in the evidence-of-depression tasks compared to annotators A_1 and B_1 . At the same time, while annotator B_2 is more inclined to indicate they are uncertain about an example than annotator A_1 , we note that annotator B_1 appears to have a higher baseline threshold of what constitutes evidence of depression than annotator A_1 . The latter is demonstrated by the fact that nearly all examples marked in the affirmative by B_1 were also marked as such by A_1 , but a large number of examples marked in the affirmative by A_1 were marked as not containing evidence of depression by B_1 .

With respect to the remission status task (bot-

¹Note that the three-class evidence-of-depression grouping simply merges the Some Evidence and Strong Evidence categories of the four-class version.

tom subplot of Figure 1), we note that annotator B_1 is more likely to mark an example as uncertain and more likely to mark an example as being inremission than annotators A_1 and B_2 . Broadly, this distribution highlights the difficulty of distinguishing active cases of clinical depression from prior experiences and lingering effects. It also serves as support for our recommendation in §6 that researchers should attempt to include the time a diagnosis was received by an individual when curating new datasets.

We acquire additional context for our results by examining the distribution of annotations as a function of the original CLPsych labels. Examining the results visualized in Figure 2, we first note that annotator A_1 classifies instances most accurately (under the assumption that ground truth is fixed over time). We believe this outcome to be a result of exposure bias; the annotation task was conducted after the completion of several modeling experiments, through which annotator A_1 was uniquely provided an opportunity to learn more about the presentation of depression by individuals in the 2015 CLPsych Shared Task dataset. We also note the distribution of "Uncertain" decisions from annotator B_2 concentrating within the original depression group. This seems to suggest annotator B_2 adopted a conservative coding approach when presented with instances that contained smaller degrees of evidence, whereas annotators A_1 and B_1 required a lower threshold of evidence to make a decision.

To conclude our reliability analysis, we examine agreement regarding the manner in which each annotator made their decision (i.e., evidence identification). We find that annotators A_1 and B_1 generally identify diagnosis disclosures within the same instances. Annotator B_2 often abstained from making a decision when presented with a disclosure due to uncertainty regarding the subject of the diagnosis. Annotator A_1 also indicated the presence of a depressed and/or irritable mood at a significantly higher rate than the other annotators, seemingly more sensitive to extreme negative emotions than the other annotators.

Discussion. Considering the difficulty of the annotation task, it is perhaps not surprising to have observed less than perfect annotator agreement. Machine learning classifiers often require hundreds of posts to make an accurate estimate of an individual's mental health status, while our annotators

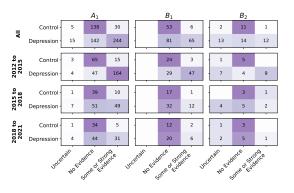


Figure 2: Distribution of annotations for the evidence of depression task (three-class) as a function of the original CLPsych labels. Affirmative evidence becomes less prevalent in the new time periods compared to the original time period for each annotator.

were only provided at maximum of 30 posts and encouraged to rely on varying levels of prior knowledge regarding the presentation of depression in social media. Critically, we emphasize that the goal of the analysis presented in §5 is *not* to curate ground truth labels of mental health status or act as clinical experts, but rather to understand biases that may exist in a depression dataset generated using self-disclosed diagnoses. The analysis of inter-rater reliability presented above provides an opportunity to further ground the results discussed in §5 and highlight areas that may benefit from future research.

A.3 Evidence Distribution

We include a breakdown of evidence annotations for individuals displaying some evidence of depression (§5) in Figure 3. Exemplary tweets for each of the evidence categories (paraphrased to maintain anonymity) are provided in Table 2. Both can be found on the following pages.

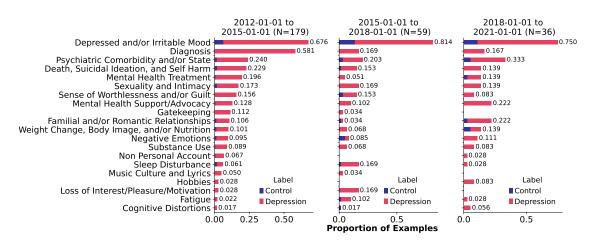


Figure 3: Distribution of evidence amongst individuals indicated as displaying at least some evidence of a depression diagnosis. Depressed and/or irritable mood is consistently the most common type of evidence within each of the three time periods.

Evidence	Exemplary Tweets
Diagnosis Disclosure	"Bipolar disorder and depression. My doctor finally agrees." "I have suffered from depression for several years now"
Depressed & Irritable Mood	"No one ever asks if I'm doing fine."
Loss of Interest/Pleasure/Motivation	"You don't understand what I'm dealing with. Get fucked." "realizing you aren't passionate about the things you used to enjoy"
Weight, Body Image, & Nutrition	"Not that anyone cares, but I'm almost at my goal weight." "I bought the dress I've always wanted, but still don't feel pretty."
Sleep Disturbance	"I CANT SLEEP. EVERYTHING HURTS. JUST LIKE ALWAYS." "Shit! Surviving on only a couple of hours of sleep again:/"
Fatigue	"mentally drained from this pandemic" "This should be effortless but I can't work any harder"
Sense of Worthlessness & Guilt	"when you let someone do anything to you" "It truly is always my fault. I probably suck."
Impaired Thought	"I'm failing my classes because I'm depressed." "at work. cant focus doe"
Death & Self Harm	"My scars are fadedunless you care to look close" "I wish you all never see a loved one fade away."
Cognitive Distortions	"Going to fail this exam. SCREWED." "I always think my bf is going to leave me"
Treatment	"Scared to tell a women that I'm in therapy" "Slowly weaning of the prozac."
Gatekeeping	"depression isn't just a bad day. fuck you all." "LET ME SHOW YOU WANT DEPRESSION IS"
Sexuality and Intimacy	"Who wants to come take some pics of me for only fans?;)" "Every girl should watch porn with their bf"
Negative Emotions	"hi sunshine! Too bad no one to spend today with." "I feel like no one cares even though I know they do"
Coping Strategies	"Have you talked to anyone about it yet?" "Art is always the easiest way to distract me from my anxiety"
Psychiatric Comorbidity & State	"Really stressing today. Lots of built up anger" "I am anorexic and cut myself."
Non-psychiatric Comorbidity	"Could use a little bit of aid #DisabilityAid" "Lots of back pain ruining what should be a beautiful day."
Substance Use	"I really shouldn't be drunk this early." "Weed makes the dreams go away and thats a good thing."
Support & Advocacy	"If I can manage a smile, I believe you can too one day!" "RIP Chester. If you're going through pain, reach out to me."
Personality and Identity	"Girls say they love a man in uniform until they do their job" "Lol grandma still think I'm bringing a boy home"
Music Culture & Lyrics	"#FallingInReverse :D" "Scene doesn't mean emo idiots. I dont want to kill myself."
Familial/Romantic Relationships	"when bae dont answer the phone xx" "Mom: You'll never lose weight. Me: Is that why dad left?"
Political & Moral Beliefs	"look in the mirror if you're not upset a cop can murder" "Trump will kill us all"
Hobbies	"Missin the old days when eveyone played Pokemon yellow" "Boys that watch the Kardashians. Love."
Non-personal Accounts	"My life was about to fall apart until I found the Calm app" "Breaking News: 5-alarm fire just outside Tulsa"

Table 2: Exemplary tweets and phrases (modified to preserve anonymity) for each of the 25 evidence categories.