# Quantifying Psycholinguistic and Mood Instability in Twitter

[1]

[1]University of Melbourne, Australia;

**Abstract**

The high prevalence of mental disorders has highlighted the importance of improving intervention strategies to early recognise the onset of mental disorders. As social media begins to emerge and people share their personal struggles online, researchers innovate in quantifying textual features to infer mental health conditions. This study aims to investigate the capabilities of textual analysis on social media to characterise people with mental disorders. The study utilised a large historical Twitter-STMHD dataset with eight disorders and a control group. We examined psycholinguistic with LIWC and temporal analysis with sentiment analysis to explore the distinguishing features. We also explored predictive analysis, achieving the best performance of 65% F1-score for binary classification and 19% in multi-class classification.

**Introduction**

Suicide as one of the leading causes of death among young people is often associated with the presence of mental disorders[1]. Approximately 970 million people globally suffered from at least one mental disorder before COVID-19. The COVID-19 pandemic further exacerbated mental health issues, particularly depression and anxiety. Despite the increase in mental health issue rates, more than 75% of individuals with mental disorders remain untreated[1]. One of the intervention problems stems from challenges in the early recognition of symptoms, including indeterminate symptoms, lack of awareness, and limited resources. There is indeed a need to improve the awareness of early symptoms for effective intervention.

Concurrently, social media enables people to share their personal experiences worldwide, including mental health struggles. Social media, with its nature as a historical, large-scale, and expressive dataset, offers valuable insights into mental health conditions. Previous research demonstrated the effectiveness of extracting textual features from social media to address mental health conditions. The research showed significant differences in linguistic style between individuals with and without mental disorders when expressing themselves online[2–6]. Further analysis revealed indications of mood instability and distinctive levels of affective in individuals with mental disorders[2,7,8]. Additionally, studies utilising social features, such as daily frequency and interactions among users, offered valuable insights[5,9]. By leveraging text processing, social media holds promise for mental health prevention and intervention efforts.

This study aims to explore the capabilities of text processing to distinguish people with and without mental disorders in social media context. We utilised the Twitter-STMHD dataset[2] and implemented the psycholinguistics analysis to infer linguistic style. We also extended the temporal analysis to explore mood instability. Finally, we combined the analysis into a set of features and inferred prediction with machine learning. However, it is important to note that the primary objective of this study was not to develop a diagnostic tool. Instead, our focus was on designing a supporting tool that could aid in characterising individuals experiencing the onset of mental disorders.

**Related Work**

Numerous studies have focused on quantifying linguistic analysis in social media to assess mental health conditions. Studies were conducted on platforms such as Reddit[3,4,10–13], Facebook[7], and Twitter[2,5,8,9]. Practically, Reddit dataset offered extensive textual analysis due to the presence of longer text. Analysing topics and linguistic style is naturally as advantageous in this platform[3,4,10]. Meanwhile, studies with Twitter datasets demonstrated the benefits of temporal analysis, such as mood instability[8] and daily activities pattern[5,8,9]. The datasets were typically large, consisting of 30K users classified into mental disorders sufferers

and control group. While mood instability was often examined in the context of borderline personality and bipolar disorder[7,8], many studies encompassed a range of other disorders.

One notable dataset was the Social Media and Mental Health Dataset (SHMD) for Reddit[4], covering conversations from 2006 to 2017. The dataset consisted of conversations from 360K users in nine subreddits. Additionally, there were other well-known datasets covering older and shorter timelines[5,10]. Lastly, SelfReported Temporally-Contextual Mental Health Diagnosis Dataset (STMHD) for Twitter[2] offered high-quality assurance through hand annotation and validation. This dataset, spanning from 2017 to 2021, captured tweets both before and after the COVID-19 pandemic.

Psycholinguistic analysis, specifically vocabulary analysis, was utilised to examine the emotional, cognitive, and structural components in individuals' writing[2–4,7,9,10]. Researchers employed tools like Linguistic Inquiry and Word Count (LIWC) to examine linguistic patterns. Regarding emotions, the control group tended to exhibit less intense expressions of anger and sadness and more prominent positive affects[2]. Other studies showed higher presence of pronouns in people with mental disorders[2,3]. Researchers also discovered that individuals with mental disorders tended to display less proficiency in grammar and syntax expressions and talked less about life-related concerns, such as money and home, than control group[2,4,7]. These findings highlight the importance of psycholinguistic analysis in understanding mental health conditions.

Temporal analysis has gained popularity in textual analysis conducted on social media. The studies ranged from mood instability[3,8,11] as well as engagement and interaction[5,9]. Researchers employed sentiment polarity as a key feature to detect mood instability in bipolar and borderline personality[8]. They aligned sentiment timelines to demonstrate the periodic mental health status. However, the research itself provided no statement about the performance of this method. Additionally, other studies demonstrated sentiment analysis to predict mental health status in non-temporal settings[8,11]. These findings underscore the potential of sentiment analysis in capturing mood patterns.

Some studies have also employed classification techniques to identify mental health conditions in social media. Traditional classification methods achieved approximately 72% F1-score in detecting depression[9]. Meanwhile, multilabel classification to detect multiple disorders reached 40-50% accuracy[4–6,10]. Advanced models such as convolutional neural networks have shown improved performance with 70% accuracy[14].

**Methods**

In this section, we outline our approach to collect and analyse data. We provide the research methodology in detail, including information on the data sources, analysis, and model design. This section ensures transparency and reproducibility of the research process.

***Data sources***

We utilised the Social Media and Mental Health Dataset (STMHD)[2]. The data was collected from Twitter and focused on eight mental disorders. The dataset consisted of 26K users identified as having one of the disorders stated in Table 1. For convenience, we named these groups as the condition groups. Additionally, there were 8K control users who were least likely to have a disorder. Figure 1 yielded there were users belonging to more than one disorder. The dataset contained user profiles and their tweets spanning from January 2017 to May 2021, covering period before and after COVID-19. The dataset was gathered through self-reported diagnosis data anchor with specific patterns, e.g. "diagnosed with <disorder name>". False anchor tweets were removed through two-step processes. First, 60% of the data underwent hand annotation and were validated by clinical psychologists. The remaining 40% of the data was filtered through pattern-matching derived from hand annotation. Dataset statistics are available in Table 1.

***Privacy and Ethic Statement***

The dataset was publicly available on the Zenodo platform[15]. The study adhered to ethical considerations to ensure individuals' privacy[2]. The dataset consisted of public, anonymised, and de-identified information. This
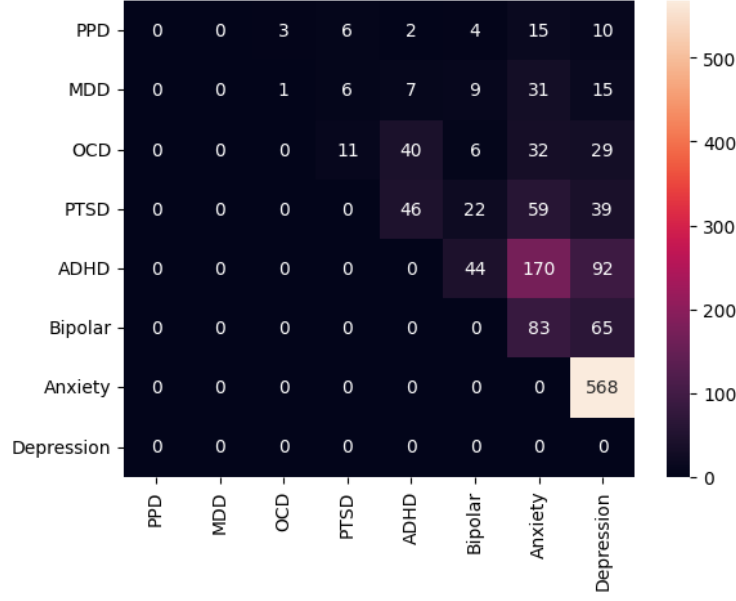
**Figure 1.** Commorbidity in dataset

**Table 1.** Statistics of dataset

| Group | #Users | Status ratio | Favourites ratio | Followers ratio | Friends ratio |
|---|---|---|---|---|---|
| Neg (control group) | 8199 | $25321 \pm 66344$ | $21763 \pm 51291$ | $818 \pm 1006$ | $844 \pm 999$ |
| ADHD (attention deficit hyperactivity disorder) | 8095 | $19060 \pm 32750$ | $34310 \pm 60030$ | $782 \pm 1009$ | $750 \pm 930$ |
| Depression | 6803 | $23177 \pm 38986$ | $35062 \pm 57404$ | $708 \pm 936$ | $774 \pm 937$ |
| Anxiety | 4843 | $20479 \pm 32222$ | $28832 \pm 50205$ | $810 \pm 1029$ | $905 \pm 1085$ |
| PTSD (post traumatic stress disorder) | 3414 | $21603 \pm 43321$ | $38871 \pm 62137$ | $648 \pm 875$ | $807 \pm 922$ |
| Bipolar disorder | 1651 | $20500 \pm 33312$ | $27995 \pm 49069$ | $757 \pm 981$ | $896 \pm 1059$ |
| OCD (obsessive compulsive disorder) | 1325 | $19180 \pm 34714$ | $33533 \pm 52270$ | $716 \pm 918$ | $800 \pm 941$ |
| MDD (major depressive disorder) | 325 | $27775 \pm 57409$ | $31266 \pm 52082$ | $732 \pm 943$ | $741 \pm 881$ |
| PPD (postpartum depression) | 247 | $21235 \pm 37713$ | $21683 \pm 44739$ | $743 \pm 974$ | $794 \pm 991$ |

study involved no direct interactions or interventions with human subjects. Therefore, obtaining informed consent was not applicable. The risk associated with this research was minimal.

### *Psycholinguistic Analysis*

We conducted a closed-vocabulary analysis to examine the linguistic features between each condition group and the control group. We utilised the well-established lexicon, LIWC (version 15)[16], which categorises vocabularies into psychological and grammatical aspects. A range of aspects were explored, including grammatical, affective, cognitive, and more. Additionally, we included categories describing daily lifestyle and personal concerns, such as money, religion, and biological processes. Table 2 summarises the categories we used. To quantify the psycholinguistic feature, we calculated the proportion of words belonging to each category aggregated by each user. Through Formula 1, we obtained the empirical distribution of users in each group. We compared the distribution in each condition group and the control group through an independent sample $t$-test.
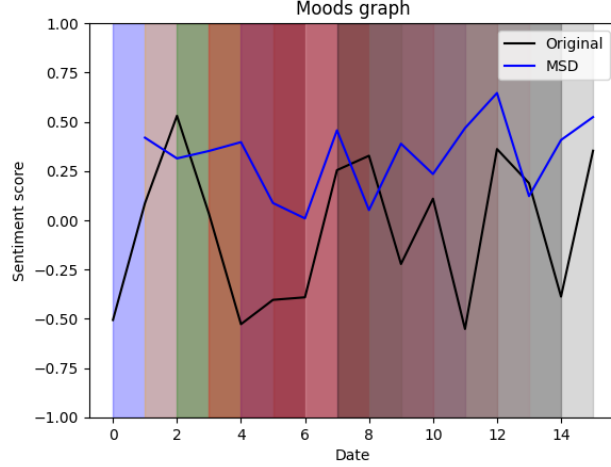
**Figure 2.** Illustration of moving standard deviation

$$F_{n_u}(x^c) = \frac{\sum_i^{n_u} 1_i^c}{n_u} \quad \forall c \in C, \forall u \in U \tag{1}$$

$$1_i^c : \text{indicator the } i\text{-th tokens is in lexicon category } c \tag{2}$$

$$n_u : \text{total number of tokens generated by user } u \tag{3}$$

### *Mood Instability*

As the moods reflect the affective state, we attempted to quantify moods with sentiment polarity analysis. In this study, we defined good moods as positive sentiments while bad moods as negative sentiments. We aimed to capture the presence of mood instability in condition groups by assuming the control group has more stable moods. Additionally, we extended the analysis to further investigate decreased mental health conditions in the control group before and after COVID-19 [17]. We performed sentiment analysis using Valence Aware Dictionary and sEntiment Reasoner (VADER) algorithm [18] which known for its effectiveness on social media.

For each tweet, we obtained a compound score of sentiment analysis. Daily moods were the average aggregation of scores for each user within a given day. We divided the time series into windows of dates with size $k = 2$. To capture instability variations, we utilised the value of moving standard deviation (MSD). For each window, we calculate its standard deviation and derived the maximum and average MSD in each user.

To further explore moods tendencies, we defined moods-ratio and moods-combo [8]. Moods-ratio, divided into positive and negative-ratio, is the proportion of days with good and bad moods respectively (See Formula 4 and 5). We also implemented the concept of moods-combo to capture the onset of hypomania/mania and depression. Moods-combo was divided into the positive-combo and negative-combo as illustrated in Figure 3. It refers to the frequency with which individuals experienced at least six consecutive days of good or bad moods. Our expectation was that condition groups would exhibit more negative tendencies and more negative-combo as indications of mental health struggles.

$$\text{positive-ratio} : F_{n_u}(x \geq 0) = \frac{\sum_i^{n_u} 1_{x_i \geq 0}}{n_u} \quad \forall c \in C, \forall u \in U \tag{4}$$

$$\text{negative-ratio} : F_{n_u}(x < 0) = \frac{\sum_i^{n_u} 1_{x_i < 0}}{n_u} \quad \forall c \in C, \forall u \in U \tag{5}$$

**Table 2.** LIWC 2015 chosen categories

| Category | Examples |
|---|---|
| Grammars | |
|    Personal pronouns | I, them, her |
|    Impersonal pronouns | it, it's, those |
|    Conjunctions | and, but, whereas |
|    Comparisons | greater, best, after |
|    Common adjectives | free, happy, long |
|    Common adverbs | very, really |
| Affective prorcesses | |
|    Positive emotion | love, nice, sweet |
|    Negative emotion | hurt, ugly, nasty |
| Cognitive processes | |
|    Insight | think, know |
|    Causation | because, effect |
|    Discrepancy | should, would |
|    Tentative | maybe, perhaps |
|    Certainty | always, never |
|    Differentiation | hasn't, but, else |
| Perceptual processes | |
|    See | view, saw, seen |
|    Hear | listen, hearing |
|    Feel | feels, touch |
| Time orientations | |
|    Past focus | ago, did, talked |
|    Present focus | today, is, now |
|    Future focus | may, will, soon |
| Biological processes | |
|    Body | cheek, hands, spit |
|    Health | clinic, flu, pill |
|    Sexual | horny, love, incest |
|    Ingestion | dish, eat, pizza |
| Social processes | |
|    Family | daughter, dad, aunt |
|    Friends | buddy, neighbor |
| Personal concerns | |
|    Work | job, majors, xerox |
|    Leisure | cook, chat, movie |
|    Home | kitchen, landlord |
|    Money | audit, cash, owe |
| Drives | |
|    Affiliation | ally, friend, social |
|    Achievement | win, success, better |
|    Power | superior, bully |
|    Reward | take, prize, benefit |
|    Risk | danger, doubt |

*Classification*

We constructed two scenarios of classification. The first scenario involved binary classification, combining condition groups into one label ($y = 1$) while the control group has the opposite ($y = 0$). Secondly,
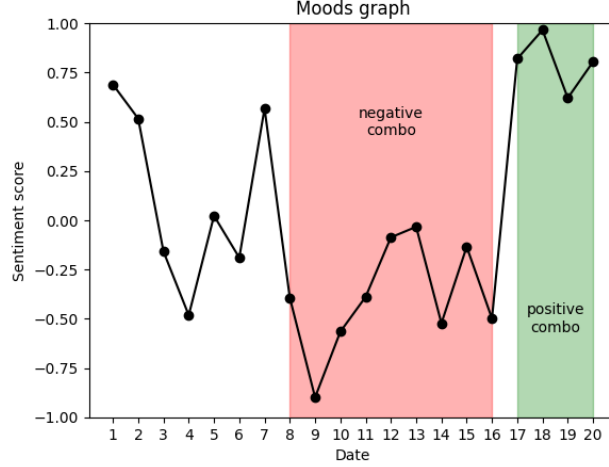
**Figure 3.** Illustration of moods-combo features

we performed multi-class classification, treating each condition group as separate labels. We utilised the calculation from previous sections as features. Additionally, we incorporated metadata aggregation, including average counts of likes, quotes, and more. We divided the data into training and testing sets in an 80:20 ratio, stratified on condition group labels. For classification, we employed traditional classifiers including *logistic regression*, *naïve bayes*, *support vector machine*, and *random forest*. Prior to the testing phase, we tuned the classifiers using cross-validation. Finally, we evaluated the performance of the classifiers using accuracy, precision, recall, and F1-score metrics (See Table 3).

**Table 3.** Confusion matrix

|  |  | *Ground-truth* | | |
|---|---|---|---|---|
|  |  | Positive | Negative | Total |
| Predicted | Positive | $TP$ | $FN$ | $TP + FN$ |
|  | Negative | $FP$ | $TN$ | $FP + TN$ |
|  | Total | $TP + FP$ | $FN + TN$ | $N$ |

1. accuracy: $\frac{TP+TN}{N}$
2. precision: $\frac{TP}{TP+FP}$
3. recall: $\frac{TP}{TP+FN}$
4. F1-score: $2\frac{precision \times recall}{precision+recall}$

## Results

In this section, we present the result of our analysis. Our analysis revealed compelling evidence for the effectiveness of our textual extraction. Statistical significance is available in details.
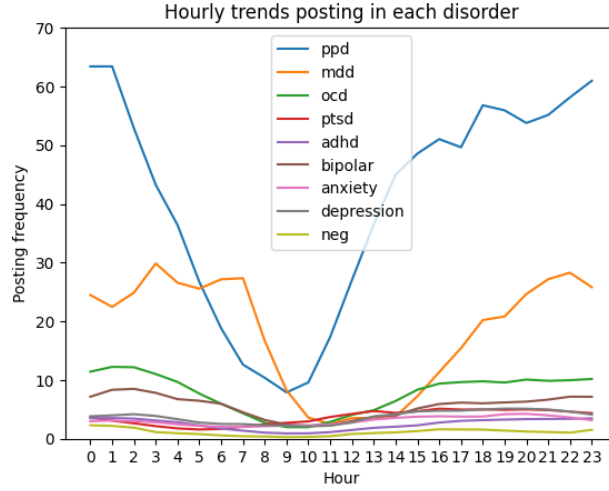
### *Features Analysis*

Initially, we examined the metadata of the dataset, as presented in Table 4. The green-coloured cell indicated higher mean values within the condition group. From the analysis, we observed significant increases in the number of replies and likes for tweets in condition groups. Conversely, the number of mentioned users decreased, with the exception of PTSD group. There were as well fewer interactions in tweets from bipolar group and more media attached from ADHD group. Lastly, we duplicated findings of hourly posting frequency

**Table 4.** Statistical significance of metadata feature

|  | PPD | MDD | OCD | PTSD | ADHD | Bipolar | Anxiety | Depression |
|---|---|---|---|---|---|---|---|---|
| Likes count | - | *** | *** | *** | *** | - | *** | *** |
| Quote count | - | * | - | - | - | ** | - | - |
| Reply count | *** | *** | *** | *** | *** | *** | *** | *** |
| Retweet count | - | * | - | - | - | . | - | - |
| Total mentioned users | - | *** | *** | *** | * | ** | *** | *** |
| Number of medias attached | - | - | - | - | *** | *** | - | - |

*** < 0.001    ** < 0.01    * < 0.05    . < 0.1



**Figure 4.** Hourly posting trends

from the original research as shown in Figure 4. It exhibited a pattern of sleep disturbance in condition groups which tended to post midnight.

The result of our psycholinguistic analysis is shown in Table 5. It indicated significant increases in most vocabulary categories for PPD, MDD, OCD, PTSD, and bipolar groups. Conversely, a decreasing trend was observed in the anxiety group for most categories. We also noticed that there were slight decreases in the vocabularies of depression groups. Surprisingly, the patterns of significance were relatively constant across categories and there were only slight differences between ADHD and control groups.

Table 6 present the analysis for mood instability features. As expected, the results yielded fluctuating moods in condition groups as evidenced by higher values of MSD. Additionally, there were significant negative tendencies in tweets made by depression, bipolar, PTSD, and ADHD groups. We also noticed a pattern of hypomania/mania in people with bipolar, anxiety, ADHD, PPD, and MDD groups. Interestingly, the positive-ratio increased across all condition groups and most of them showed fewer signs of negative-combo (depression signs), contrary to our initial expectations. Additionally, we discovered that tweets posted after the COVID-19 period exhibited a lower positive tendency.

### *Classification*

Table 7 and 8 yielded the performance of our classification models. As past papers indicated, the performance of binary classification was higher than multi-class classification. In terms of binary classification, the performance yielded 70-80% accuracy and 45-65% F1-score. Meanwhile, multi-class performed at 25-39% accuracy and 15-19% F1-score. Overall, *random forest* stood out with the best performance of 65% for binary

**Table 5.** Statistical significance of psycholinguistic analysis

| | PPD | MDD | OCD | PTSD | ADHD | Bipolar | Anxiety | Depression |
|---|---|---|---|---|---|---|---|---|
| **Grammars** | | | | | | | | |
| Personal pronouns | *** | *** | *** | *** | - | *** | ** | - |
|   i | *** | *** | *** | *** | * | *** | - | - |
|   we | *** | *** | *** | * | - | - | ** | * |
|   you | *** | *** | *** | *** | - | *** | *** | - |
|   shehe | *** | *** | *** | ** | - | . | ** | - |
|   they | *** | *** | *** | *** | - | * | ** | - |
| Impersonal pronouns | *** | *** | *** | *** | - | ** | ** | - |
| Common adverbs | *** | *** | *** | *** | - | *** | ** | - |
| Conjunctions | *** | *** | *** | *** | - | *** | * | - |
| Common Adjectives | *** | *** | *** | *** | - | ** | *** | - |
| Comparisons | *** | *** | *** | *** | - | * | *** | - |
| **Affective processes** | *** | *** | *** | *** | - | *** | ** | - |
| Positive emotion | *** | *** | *** | *** | - | *** | *** | - |
| Negative emotion | *** | *** | *** | *** | - | ** | *** | - |
|   Anxiety | *** | *** | *** | *** | - | *** | * | - |
|   Anger | *** | *** | *** | *** | - | ** | * | - |
|   Sad | *** | *** | *** | *** | - | *** | * | - |
| **Social processes** | *** | *** | *** | *** | - | ** | *** | - |
| Family | *** | *** | *** | *** | - | ** | *** | - |
| Friends | *** | *** | *** | *** | - | ** | *** | - |
| **Cognitive processes** | *** | *** | *** | *** | - | ** | ** | - |
| Insight | *** | *** | *** | *** | - | ** | ** | - |
| Causation | *** | *** | *** | *** | - | ** | ** | - |
| Discrepancy | *** | *** | *** | *** | - | ** | ** | - |
| Tentative | *** | *** | *** | *** | - | ** | ** | - |
| Certainty | *** | *** | *** | *** | - | ** | *** | - |
| Differentiation | *** | *** | *** | *** | - | ** | ** | - |
| **Perceptual processes** | *** | *** | *** | *** | - | ** | *** | - |
| See | *** | *** | *** | - | - | * | *** | * |
| Hear | *** | *** | *** | *** | - | *** | ** | - |
| Feel | *** | *** | *** | *** | - | *** | * | - |
| **Biological processes** | *** | *** | *** | *** | - | *** | * | - |
| Body | *** | *** | *** | *** | - | *** | ** | - |
| Health | *** | *** | *** | *** | - | *** | - | - |
| Sexual | *** | *** | *** | *** | * | *** | - | . |
| Ingestion | *** | *** | *** | - | - | - | * | - |
| **Drives** | *** | *** | *** | ** | - | * | *** | . |
| Affiliation | *** | *** | *** | *** | - | * | *** | . |
| Achievement | *** | *** | *** | * | - | - | *** | * |
| Power | *** | *** | *** | ** | - | - | *** | . |
| Reward | *** | *** | *** | * | - | . | *** | * |
| Risk | *** | *** | *** | *** | - | * | ** | - |
| **Time orientation** | | | | | | | | |
| Past focus | *** | *** | *** | *** | - | ** | ** | - |
| Present focus | *** | *** | *** | *** | - | ** | *** | - |
| Future focus | *** | *** | *** | *** | - | ** | *** | - |
| **Personal concerns** | | | | | | | | |
| Work | *** | *** | *** | * | - | - | ** | * |
| Leisure | *** | *** | *** | | - | - | *** | ** |
| Home | *** | *** | *** | . | - | - | ** | - |
| Money | *** | *** | *** | - | . | - | ** | . |
| Religion | *** | *** | *** | * | - | ** | * | - |
| Death | *** | *** | *** | * | - | - | *** | - |

**Table 6.** Statistical significance of mood instability feature

| | PPD | MDD | OCD | PTSD | ADHD | Bipolar | Anxiety | Depression | After covid |
|---|---|---|---|---|---|---|---|---|---|
| Max MSD | *** | *** | ** | *** | *** | *** | *** | *** | - |
| Average MSD | *** | * | *** | *** | *** | *** | *** | *** | - |
| Positive ratio | *** | * | *** | *** | *** | *** | *** | *** | *** |
| Negative ratio | - | - | - | *** | *** | * | - | *** | - |
| Positive combo | *** | *** | - | . | *** | *** | *** | - | *** |
| Negative combo | * | - | *** | *** | *** | - | *** | *** | *** |

*** $< 0.001$     ** $< 0.01$     * $< 0.05$     . $< 0.1$

classification and 19% for multi-class classification. In multi-class scenario, we discovered that our model tended to classify individuals into PPD and Bipolar groups as evidenced by the higher predicted values in both true positives and false positives (See Figure 5).

**Table 7.** Binary classification result

| | Acc | P | R | F1 |
|---|---|---|---|---|
| Logistic Regression | 0.77 | 0.67 | 0.51 | 0.46 |
| Naive Bayes | 0.72 | 0.55 | 0.53 | 0.53 |
| Support Vector Machine | 0.77 | 0.64 | 0.51 | 0.45 |
| Random Forest | 0.80 | 0.75 | 0.63 | 0.65 |

**Table 8.** Multi-class classification result

| | Acc | P | R | F1 |
|---|---|---|---|---|
| Logistic Regression | 0.35 | 0.29 | 0.19 | 0.17 |
| Naive Bayes | 0.25 | 0.2 | 0.18 | 0.15 |
| Support Vector Machine | 0.36 | 0.26 | 0.2 | 0.18 |
| Random Forest | 0.39 | 0.26 | 0.21 | 0.19 |

**Discussion**

In most condition groups' interaction patterns, we noticed higher numbers of replies and likes, coupled with lower numbers of mentioned users. We suggested the tweets might generate significant engagement among users without direct interactions with the specific user. However, the finding was aligned with previous studies stating fewer interactions among users[5,9]. Conversely, the higher numbers of mentioned users and lower counts of followers in PTSD tweets might indicate intense engagement. It is possibly associated with specific attachment styles seen in individuals with PTSD, particularly those with Complex PTSD (C-PTSD) which have not been explored before.

For mood instability features, the depression-related groups exhibited decreasing trend in negative-combo values with higher negative-ratio. This suggested long exposure to depressive episodes with less frequent mood-shifting and denser episodes. This finding indicated the persistence and chronicity of depression in certain individuals. As for higher positive-ratio in condition groups, we found this insight interesting. We supposed individuals might present a curated version of themselves, including portraying positive images or as a way of coping mechanism. However, further sampling might be needed to explore the patterns. Additionally, the analysis of post COVID-19 tweets suggested there was an indication of worse mood stability in the control group as shown in previous studies[17]. Overall, the features extend the result of previous research[8], proving that the features showed valuable indicators.
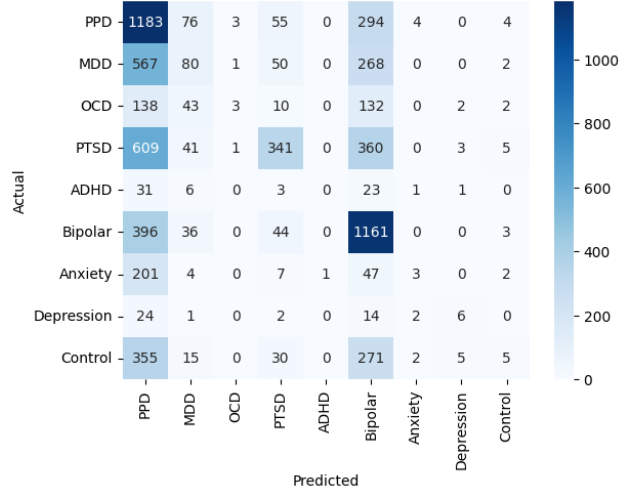
**Figure 5.** Confusion matrix of multi-class classification

Our analysis of psycholinguistics features yielded anomalous results. Further analysis indicated a limitation in our study, as unforeseen flaws in the code were discovered due to parallel data processing on ADHD, depression, and control groups. However, we could not revise prior to submitting the paper due to time constraints. To further validate the features and provide better results, we experimented classifications without psycholinguistic features. Table 9 and 10 show the result of the experiment. We discovered that the performance provided no significant differences. We supposed the psycholinguistic features gave slight contributions to the performance due to flaws.

**Table 9.** Binary classification result revised

|                          | Acc  | P    | R    | F1   |
|--------------------------|------|------|------|------|
| Logistic Regression      | 0.77 | 0.68 | 0.51 | 0.46 |
| Naive Bayes              | 0.72 | 0.55 | 0.53 | 0.53 |
| Support Vector Machine   | 0.77 | 0.66 | 0.51 | 0.46 |
| Random Forest            | 0.80 | 0.75 | 0.64 | 0.66 |

**Table 10.** Multi-class classification result revised

|                          | Acc  | P    | R    | F1   |
|--------------------------|------|------|------|------|
| Logistic Regression      | 0.35 | 0.18 | 0.18 | 0.15 |
| Naive Bayes              | 0.24 | 0.15 | 0.14 | 0.13 |
| Support Vector Machine   | 0.36 | 0.17 | 0.18 | 0.15 |
| Random Forest            | 0.40 | 0.21 | 0.2  | 0.17 |

**Conclusion**

Overall, we examined text processing techniques on the Twitter-STHMD dataset to distinguish individuals with mental disorders and control groups. We performed psycholinguistic and temporal analysis which were combined for predictive machine learning. Our predictive model employing random forest achieved the highest F1-score of 65% for binary classification and 19% for multi-class.

In future work, some areas warrant further exploration. Firstly, it is worth revisiting the psycholinguistic

analysis and fixing associated flaws. Secondly, the prominence of time series moods would be valuable to investigate the aligning time series in each group. This would provide insights into temporal dynamics and patterns and further enable early recognition of symptoms based on mood fluctuations. Additionally, incorporating POS tagging and syntactic analysis tools might offer a more detailed and quantitative understanding of linguistic analysis. Furthermore, topical modeling techniques could be employed to uncover latent themes and topics within the dataset.

**References**

1. Organization WH. World mental health report: transforming mental health for all. World Health Organization; 2022.

2. Suhavi, Singh AK, Arora U, Shrivastava S, Singh A, Shah RR, et al. Twitter-STMHD: An Extensive User-Level Database of Multiple Mental Health Disorders. Proceedings of the International AAAI Conference on Web and Social Media. 2022 May;16:1182-91. Available from: https://ojs.aaai.org/index.php/ICWSM/article/view/19368.

3. Gkotsis G, Oellrich A, Hubbard T, Dobson R, Liakata M, Velupillai S, et al. The language of mental health problems in social media. In: Proceedings of the Third Workshop on Computational Lingusitics and Clinical Psychology. San Diego, CA, USA: Association for Computational Linguistics; 2016. p. 63-73. Available from: http://aclweb.org/anthology/W16-0307.

4. Cohan A, Desmet B, Yates A, Soldaini L, MacAvaney S, Goharian N. SMHD: A Large-Scale Resource for Exploring Online Language Usage for Multiple Mental Health Conditions. arXiv; 2018. ArXiv:1806.05258 [cs]. Available from: http://arxiv.org/abs/1806.05258.

5. 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. Baltimore, Maryland, USA: Association for Computational Linguistics; 2014. p. 51-60. Available from: http://aclweb.org/anthology/W14-3207.

6. Gaur M, Kursuncu U, Alambo A, Sheth A, Daniulaityte R, Thirunarayan K, et al. "Let Me Tell You About Your Mental Health!": Contextualized Classification of Reddit Posts to DSM-5 for Web-based Intervention. In: Proceedings of the 27th ACM International Conference on Information and Knowledge Management. Torino Italy: ACM; 2018. p. 753-62. Available from: https://dl.acm.org/doi/10.1145/3269206.3271732.

7. Saha K, Chan L, De Barbaro K, Abowd GD, De Choudhury M. Inferring Mood Instability on Social Media by Leveraging Ecological Momentary Assessments. Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies. 2017 Sep;1(3):1-27. Available from: https://dl.acm.org/doi/10.1145/3130960.

8. Saravia E, Chang CH, De Lorenzo RJ, Chen YS. MIDAS: Mental illness detection and analysis via social media. In: 2016 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM). San Francisco, CA, USA: IEEE; 2016. p. 1418-21. Available from: http://ieeexplore.ieee.org/document/7752434/.

9. De Choudhury M, Gamon M, Counts S, Horvitz E. Predicting Depression via Social Media. Proceedings of the International AAAI Conference on Web and Social Media. 2021 Aug;7(1):128-37. Available from: https://ojs.aaai.org/index.php/ICWSM/article/view/14432.

10. Coppersmith G, Dredze M, Harman C, Hollingshead K. From ADHD to SAD: Analyzing the Language of Mental Health on Twitter through Self-Reported Diagnoses. In: Proceedings of the 2nd Workshop on Computational Linguistics and Clinical Psychology: From Linguistic Signal to Clinical Reality. Denver, Colorado: Association for Computational Linguistics; 2015. p. 1-10. Available from: http://aclweb.org/anthology/W15-1201.

11. Song H, You J. Feature Attention Network: Interpretable Depression Detection from Social Media. Information and Computation. 2018.

12. 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. San Jose California USA: ACM; 2016. p. 2098-110. Available from: https://dl.acm.org/doi/10.1145/2858036.2858207.

13. Gkotsis G, Oellrich A, Velupillai S, Liakata M, Hubbard TJP, Dobson RJB, et al. Characterisation of mental health conditions in social media using Informed Deep Learning. Scientific Reports. 2017 Mar;7(1):45141. Available from: https://www.nature.com/articles/srep45141.

14. Ive J, Gkotsis G, Dutta R, Stewart R, Velupillai S. Hierarchical neural model with attention mechanisms for the classification of social media text related to mental health. In: Proceedings of the Fifth Workshop on Computational Linguistics and Clinical Psychology: From Keyboard to Clinic. New Orleans, LA: Association for Computational Linguistics; 2018. p. 69-77. Available from: http://aclweb.org/anthology/W18-0607.

15. Suhavi, Singh A, Arora U, Shrivastava S, Singh A, Shah RR, et al.. Twitter Self-reported Temporally-contextual Mental Health Diagnosis (Twitter-STMHD) Dataset. Zenodo; 2022. Available from: https://doi.org/10.5281/zenodo.5854911.

16. Pennebaker J, Boyd R, Jordan K, Blackburn K. The development and psychometric properties of LIWC2015. University of Texas at Austin; 2015.

17. Xie Y, Xu E, Al-Aly Z. Risks of mental health outcomes in people with covid-19: cohort study. BMJ. 2022;376. Available from: https://www.bmj.com/content/376/bmj-2021-068993.

18. Hutto C, Gilbert E. VADER: A Parsimonious Rule-Based Model for Sentiment Analysis of Social Media Text. Proceedings of the International AAAI Conference on Web and Social Media. 2014 May;8(1):216-25. Available from: https://ojs.aaai.org/index.php/ICWSM/article/view/14550.

## Appendices

### Code: Transforming data from JSON to CSV

```python
import os
import sys
import json
import zipfile
import pandas as pd

# Get filename
file_arg = sys.argv[1]
dis = file_arg.split("/")[-1][:-4]
user_col = ["id","creation_timestamp","description","favourites_count","followers_count",
            "friends_count", "geo_tag","banner_link","display_image_link","status_count",
            "verified_check", "anchor_tweet", "anchor_tweet_date","disorder"]
tweet_col = ["disorder_flag", "text", "conversation_id", "tweet_id", "language",
             "likes_count", "quote_count", "reply_count", "retweet_count", "source_name",
             "timestamp_tweet", "mentionedUsers", "media", "user_id", "covid"]

# Database retrieve
if os.path.exists(f'users_{dis}.csv'):
    users_data = pd.read_csv(f'users_{dis}.csv', dtype=str, lineterminator=";")
    users = users_data[users_data.disorder == dis]['id'].values
else:
    users_data = pd.DataFrame(columns=user_col)
    users = []
print("Numbers of transformed users:", len(users))

# Filter directories (only the one which haven't processed)
covid_date = '2020-03-11'
archive = zipfile.ZipFile(f'{file_arg}', 'r')
listdirs = archive.namelist()
dirs = []
for s in listdirs:
    s_split = s.split("/")
    if len(s_split) == int(sys.argv[2]) and s.endswith("/") and s_split[-2] \
        not in users: dirs.append(s)
print("Total users left:", len(dirs))

# Function to save file, concat data only at the end
def save_file(users, tweets):
    users_data = pd.concat(users)
    tweets_data = pd.concat(tweets)
    print(len(tweets_data['user_id'][tweets_data['user_id'].isnull()]))

    if os.path.exists(f'users_{dis}.csv'):
        users_data.to_csv(f'users_{dis}.csv', mode='a', header=False,
                          index=False, lineterminator=";")
        tweets_data.to_csv(f'tweets_{dis}.csv', mode='a', header=False,
                           index=False, lineterminator=";")
    else:
        users_data.to_csv(f'users_{dis}.csv', index=False, lineterminator=";")
        tweets_data.to_csv(f'tweets_{dis}.csv', index=False, lineterminator=";")


# Run transformation
total_user = 0
while total_user != len(dirs):
    i = 0
    users_new = []
    tweets_new = []
    try:
        for dir in dirs[total_user:]:
            namespace = dir.split("/")
            user_id = namespace[-2]
            i += 1
```

```
64                total_user += 1
65                print("=", end="", flush=True)
66
67                if dis == 'neg':
68                    anchor = {}
69                else:
70                    anchor = json.load(archive.open(dir + 'anchor_tweet.json'))
71
72                user = json.load(archive.open(dir + 'user.json'))
73                tweets = json.load(archive.open(dir + 'tweets.json'))
74
75                user_data = {**user, **anchor}
76                user_df = pd.DataFrame(user_data, index=[0])
77                user_df['disorder'] = namespace[5]
78
79                tweet_temp = []
80                for keydate, tweet in tweets.items():
81                    tweet_df = pd.DataFrame(tweet)
82                    tweet_df['user_id'] = user_id
83                    tweet_df['covid'] = 1 if keydate > covid_date else 0
84                    tweet_df['media'] = tweet_df['media'] \
85                        .apply(lambda x: str([m["type"] for m in x]))
86                    tweet_df = tweet_df.reindex(columns=tweet_col).fillna("")
87                    tweet_temp.append(tweet_df)
88
89                user_df = user_df.reindex(columns=user_col).fillna("")
90                users_new.append(user_df)
91                tweets_new += tweet_temp
92
93                if i == 100:
94                    print()
95                    print("Taking a break...")
96                    break
97
98            save_file(users_new, tweets_new)
99            print("Let's do it again!", total_user, "out of", len(dirs))
100
101    # Enabling intermittent transformation with keyboard interruption
102    except KeyboardInterrupt:
103        print("Interrupt execution...")
104        save_file(users_new, tweets_new)
105        break
```

## Code: Scratches

```
1  # -*- coding: utf-8 -*-
2  """Scratch.ipynb
3
4  Automatically generated by Colaboratory.
5
6  Original file is located at
7      https://colab.research.google.com/drive/12M8xV3f0uf01sASpLj1JPVd2aLurkRs6
8  """
9
10 import matplotlib.pyplot as plt
11 import random
12
13 # Generate 20 random x-values
14 x = range(1, 21)
15
16 # Generate random y-values between -1 and 1
17 random.seed(0)
18 y = [random.uniform(-1, 1) for _ in range(20)]
19
```

```python
# Ensure at least 4 consecutive positive values
positive_start = random.randint(0, 17)
for i in range(positive_start, positive_start + 4):
    y[i] = random.uniform(0, 1)

# Ensure at least 7 consecutive negative values
negative_start = random.randint(0, 13)
for i in range(negative_start, negative_start + 7):
    y[i] = random.uniform(-1, 0)

# Plot the graph
plt.title('Moods graph')
plt.plot(x, y, marker='o', linestyle='-', color='black')

# Set y-axis limits between -1 and 1
plt.ylim(-1, 1)
plt.xticks(x)

# Add labels
plt.xlabel('Date')
plt.ylabel('Sentiment score')

plt.axvspan(8, 16, color='red', alpha=0.3)
plt.axvspan(17, 20, color='green', alpha=0.3)
plt.text(12, 0.5, 'negative\ncombo', color='black', ha='center',
         va='center', fontsize=10)
plt.text(18.5, -0.5, 'positive\ncombo', color='black', ha='center',
         va='center', fontsize=10)

# Show the plot
plt.show()

import matplotlib.pyplot as plt
import numpy as np
import pandas as pd

# Generate x values from 0 to 100
x = np.arange(0, 16)

# Generate random y values between -1 and 1
y = np.random.uniform(-0.6, 0.6, size=len(x))

# Set the sliding window size

# Define the colors for the sliding window
colors = ['blue', 'orange', 'green', 'red', 'purple', 'brown', 'pink', 'black', 'gray']

# Plot the original data
plt.plot(x, y, color='black', label='Original')

# Plot the sliding window averages with different colors
rolling_std = pd.Series(y).rolling(window=2).std()
plt.plot(x, y, color='black', label='Original')

# Set the y-axis limit
plt.ylim(-1, 1)

# Add legend and labels
plt.title('Moods graph')
plt.xlabel('Date')
plt.ylabel('Sentiment score')
k = 7
for i in range(len(colors)):
    plt.axvspan(i, i+k, color=colors[i], alpha=0.3)

# Show the plot
plt.show()
```

## Code: Analysing user profiles

```
1  # -*- coding: utf-8 -*-
2  """User Analysis.ipynb
3
4  Automatically generated by Colaboratory.
5
6  Original file is located at
7      https://colab.research.google.com/drive/18nDqWoHrTeIQ0EgApw3RE7kzczfOXm9m
8  """
9
10 DIR = "/content/drive/MyDrive/College (Master)/Semester 3/ \
11      COMP90090 - Natural Language Processing for Health/Assignment/Assignment 3/"
12
13 from google.colab import drive
14 drive.mount('/content/drive')
15
16 import pandas as pd
17
18 DISORDERS = ['ppd', 'mdd', 'ocd', 'ptsd', 'adhd',
19              'bipolar', 'anxiety', 'depression']
20 LABEL = DISORDERS + ['neg']
21
22 users = []
23
24 for dis in LABEL:
25     user = pd.read_csv(f"{DIR}/Dataset/users/users_{dis}.csv",
26                        dtype=str, lineterminator=";")
27     users.append(user)
28
29 df = pd.concat(users).reset_index(drop=True)
30 df = df[~df.disorder.isnull()]
31 df['favourites_count'] = df['favourites_count'].astype(int)
32 df['followers_count'] = df['followers_count'].astype(int)
33 df['friends_count'] = df['friends_count'].astype(int)
34 df.head()
35
36 """# Disorder counts"""
37
38 # Statistics per user
39 df['disorder'].value_counts()
40
41 combined_df = df[['id', 'disorder']].groupby('id')['disorder'] \
42                   .apply(lambda x: ','.join(map(str, x))).reset_index()
43 commorbid = combined_df[combined_df['disorder'].str.contains(',')]['disorder']
44 commorbid = commorbid.value_counts().reset_index(name="Frequency")
45 commorbid.head()
46
47 for i, disorder in enumerate(commorbid['index']):
48     dsplit = disorder.split(",")
49     if len(dsplit) > 2:
50         for j in range(len(dsplit)):
51             if j == (len(dsplit) - 1): break
52             for d in dsplit[j+1:]:
53                 comm_id = commorbid[commorbid['index'] == f"{dsplit[j]},{d}"].index
54                 commorbid.loc[comm_id, 'Frequency'] += commorbid.loc[i, 'Frequency']
55         commorbid.drop(i, inplace=True)
56
57 commorbid.head()
58
59 import matplotlib.pyplot as plt
60 import seaborn as sns
61
62 commorbid['Condition1'], commorbid['Condition2'] = commorbid['index'].str.split(',', 1).str
63 pivot_table = commorbid.pivot_table(index='Condition1', columns='Condition2',
```

```python
                                          values='Frequency', fill_value=0)
pivot_table = pivot_table.reindex(index=DISORDERS, columns=DISORDERS) \
                    .fillna(0).astype("int64")
ax = sns.heatmap(pivot_table, annot=True, fmt='g')

ax.set_xticklabels(["PPD", "MDD", "OCD", "PTSD", "ADHD",
                    "Bipolar", "Anxiety", "Depression"])
ax.set_yticklabels(["PPD", "MDD", "OCD", "PTSD", "ADHD",
                    "Bipolar", "Anxiety", "Depression"])
ax.set_xlabel(None)
ax.set_ylabel(None)
plt.show()

"""# Favorites, Followers, Friends"""

favourites = []
followers = []
friends = []
status = []

for l in LABEL:
    status.append(df[df.disorder == l].status_count.astype(int))
    favourites.append(df[df.disorder == l].favourites_count.astype(int))
    followers.append(df[df.disorder == l].followers_count.astype(int))
    friends.append(df[df.disorder == l].friends_count.astype(int))

stat = []

for i, l in enumerate(['neg', 'adhd', 'depression', 'anxiety', 'ptsd',
                        'bipolar', 'ocd', 'mdd', 'ppd']):
    print("==================== Statistics of group", l)
    stat.append([])
    print("Status count", round(status[i].mean()), "+-", round(status[i].std()))
    print("Favorites", round(favourites[i].mean()), "+-", round(favourites[i].std()))
    print("Followers", round(followers[i].mean()), "+-", round(followers[i].std()))
    print("Friends", round(friends[i].mean()), "+-", round(friends[i].std()))
```

**Code: Extracting metadata analysis**

```python
# -*- coding: utf-8 -*-
"""Metadata Analysis.ipynb

Automatically generated by Colaboratory.

Original file is located at
    https://colab.research.google.com/drive/1Q8glz1mULTdqhp40OeGup1jz4fZn-1b2
"""

DIR = "/content/drive/MyDrive/College (Master)/Semester 3/ \
        COMP90090 - Natural Language Processing for Health/Assignment/Assignment 3/"

import os
import nltk
import pandas as pd
import scipy.stats as stats
from pandas.errors import ParserError
from nltk.sentiment.vader import SentimentIntensityAnalyzer

import warnings
warnings.filterwarnings("ignore")

nltk.download('vader_lexicon')

DISORDERS = ['ppd', 'mdd', 'ocd', 'ptsd', 'adhd',
```

```
26                    'bipolar', 'anxiety', 'depression']
27  LABEL = DISORDERS + ['neg']
28
29  """# Feature Extraction"""
30
31  # Store data about user_id, tweet_id, date, sentiment, disorder, covid
32  for l in LABEL[7:]:
33      print(f"Working on label {l}")
34      reader = pd.read_csv(f"{DIR}/Dataset/tweets/tweets_{l}.csv", dtype=str,
35                           chunksize=100000, on_bad_lines='warn', lineterminator=";")
36
37      i = 0
38      for r in reader:
39          i += 1
40          df = r[['tweet_id', 'user_id', 'covid', 'disorder_flag',
41                  'likes_count', 'quote_count', 'reply_count', 'retweet_count']]
42          df['mentioned_users'] = r['mentionedUsers']
43                                  .apply(lambda x:
44                                         0 if x == "[]" or
45                                              x == ""
46                                         else len(x.split(",")))
47          df['medias'] = r['media']
48                         .apply(lambda x:
49                                0 if x == "[]" or
50                                     x == ""
51                                else len(x.split(",")))
52          print(f"Chunk number {i}")
53
54          if not os.path.exists(f"{DIR}/Dataset/MA/{l}.csv"):
55              df.to_csv(f"{DIR}/Dataset/MA/{l}.csv")
56          else:
57              df.to_csv(f"{DIR}/Dataset/MA/{l}.csv", mode='a', header=False)
58
59      print(" FINISHED")
60
61  """# Analisis"""
62
63  columns = ['likes_count', 'quote_count', 'reply_count',
64             'retweet_count', 'mentioned_users', 'medias']
65
66  dfs = []
67  hourly = []
68  types = dict(zip(columns,[int for i in range(len(columns))]))
69
70  for l in LABEL:
71      print(f"Working on label {l}")
72
73      reader = pd.read_csv(f"{DIR}/Dataset/MA/{l}.csv",
74                           dtype=types, chunksize=1000000,
75                           on_bad_lines='warn', index_col=0)
76
77      temp = [[] for i in range(len(columns))]
78
79      i = 0
80      append_df = None
81      for r in reader:
82          i += 1
83          # Moving chunks
84          if i > 1: r = pd.concat([append_df, r], ignore_index=True, axis=0)
85
86          # Get dataframe without last user
87          last_user = r.user_id.unique()[-1]
88          mask = (r.user_id == last_user)
89          append_df = r[mask]
90          r = r[~mask]
91
92          # Aggregate values
93          aggs = dict(zip(columns,['mean' for i in range(len(columns))]))
```

```python
            result = r.groupby(['user_id']).agg(aggs)
            for j in range(len(columns)):
                temp[j] += result[columns[j]].tolist()

            print("Chunk number", i)

    df = pd.DataFrame([list(i) for i in zip(*temp)], columns=columns)
    dfs.append(df)
    print("FINISHED")

summ = pd.DataFrame(index=columns, columns=DISORDERS)

for j, col in enumerate(columns):
    for i in range(len(DISORDERS)):
        t_statistic1, p_value1 = stats.ttest_ind(dfs[i][col].values,
                                                 dfs[8][col].values)
        value = "***" if p_value1 < 0.001 else "**"
                    if p_value1 < 0.01 else "*"
                    if p_value1 <0.05 else "."
                    if p_value1 < 0.1 else ""
        color = "r" if round(t_statistic1,10) < 0 else "g"
                 if round(t_statistic1,10) > 0 else ""
        summ.at[col, DISORDERS[i]] = value + color

summ

summ.to_latex()

for i, df in enumerate(dfs):
    df.to_csv(f"{DIR}/Dataset/MA/{LABEL[i]}_extracted.csv")
```

## Code: Extracting mood instability

```python
# -*- coding: utf-8 -*-
"""Mood Instability.ipynb

Automatically generated by Colaboratory.

Original file is located at
    https://colab.research.google.com/drive/1KO9pnqF3Qy4pM-kl6ITtBzYmNJ3z-jwQ
"""

DIR = "/content/drive/MyDrive/College (Master)/Semester 3/ \
        COMP90090 - Natural Language Processing for Health/Assignment/Assignment 3/"

import nltk
import numpy as np
import pandas as pd
import scipy.stats as stats
import statsmodels.api as sm
import matplotlib.pyplot as plt
import matplotlib.dates as mdates

from itertools import groupby
from pandas.errors import ParserError
from nltk.sentiment.vader import SentimentIntensityAnalyzer


import warnings
warnings.filterwarnings("ignore")

nltk.download('vader_lexicon')

DISORDERS = ['ppd', 'mdd', 'ocd', 'ptsd', 'adhd',
```

```python
32                  'bipolar', 'anxiety', 'depression']
33  LABEL = DISORDERS + ['neg']
34
35  """# Feature Extraction"""
36
37  import os
38
39  # Store data about user_id, tweet_id, date, sentiment, disorder, covid
40  sid = SentimentIntensityAnalyzer()
41
42  for l in LABEL[8:]:
43      print(f"Working on label {l}")
44      reader = pd.read_csv(f"{DIR}/Dataset/tweets/tweets_{l}.csv", dtype=str,
45                           chunksize=500000, on_bad_lines='warn', lineterminator=";")
46
47      i = 0
48      for r in reader:
49          i += 1
50          sentiments = []
51          for text in r['text'].values:
52              try:
53                  sentiments.append(sid.polarity_scores(text))
54              except AttributeError:
55                  sentiments.append({'neg':0, 'neu':0, 'pos':0, 'compound':0})
56
57          sent_df = pd.DataFrame(sentiments)
58          df = r[['tweet_id', 'user_id', 'disorder_flag', 'timestamp_tweet']]
59          df.loc[:,'neg'] = sent_df['neg'].values
60          df.loc[:,'neu'] = sent_df['neu'].values
61          df.loc[:,'pos'] = sent_df['pos'].values
62          df.loc[:,'compound'] = sent_df['compound'].values
63          print(f"Chunk number {i}")
64
65          if not os.path.exists(f"{DIR}/Dataset/NEW/{l}.csv"):
66              print("masuk")
67              df.to_csv(f"{DIR}/Dataset/NEW/{l}.csv")
68          else:
69              df.to_csv(f"{DIR}/Dataset/NEW/{l}.csv", mode='a', header=False)
70
71
72      print(" FINISHED")
73
74  """# Exploratory"""
75
76  mi = pd.read_csv(DIR + 'Dataset/MI/ppd.csv', index_col=0)
77  mi['date'] = mi['timestamp_tweet'].str.split(" ").str[0]
78  mi.drop(['timestamp_tweet'], axis=1, inplace=True)
79  print(mi.shape)
80
81  mi['date'] = pd.to_datetime(mi['date'])
82
83  result_date = mi.groupby(['user_id', 'date']).agg({'compound': 'mean',
84                                                     'neg': 'mean',
85                                                     'neu': 'mean',
86                                                     'pos': 'mean'})
87
88  print(result_date)
89
90  plt.figure(figsize=(15, 6))
91  plt.xlabel('Date')
92  plt.ylabel('Value')
93  plt.title('Time Series Data with Smoothing')
94  plt.ylim(-1,1)
95  test = []
96
97  # Define the start and end dates of the period to highlight
98  highlight_start = '2020-03-01'  # Example start date
99
```

```python
100  # Convert the highlight dates to matplotlib date format
101  highlight_start = mdates.datestr2num(highlight_start)
102
103  # Draw a rectangle to highlight the period
104
105  for user in mi['user_id'].unique()[:5]:
106      comp = result_date.loc[user, 'compound']
107      rolling_avg = comp.rolling(window=10).mean()
108      rolling_std = comp.rolling(window=10).std()
109      test.append(comp.values[~np.isnan(comp.values)])
110
111      # plt.plot(comp.index, comp.values, color='blue')
112      # plt.plot(comp.index, rolling_avg, color='green')
113      plt.plot(comp.index, rolling_std)
114
115  plt.axvspan(highlight_start, plt.xlim()[1], facecolor='red', alpha=0.3, zorder=0)
116  # ccf = sm.tsa.stattools.ccf(test[0], test[1], adjusted=False)
117  # plt.plot(comp.index, ccf, label='Rolling Average (Window = 7)')
118  plt.show()
119
120  """# Detailed Feature Extraction"""
121
122  features = ['max_msd', 'mean_msd',
123              'positive_ratio', 'negative_ratio',
124              'positive_combo', 'negative_combo']
125
126  columns = ['compound']
127
128  dfs = []
129  hourly_rate = np.zeros((len(LABEL), 24))
130
131  before_covid = []
132  after_covid = []
133  covid_date = pd.to_datetime('2020-03-11')
134
135  def extract_feature(compounds, temp):
136      # Moving windows
137      rolling_std = compounds.rolling(window=2).std()
138      temp[0].append(np.max(rolling_std))
139      temp[1].append(np.mean(rolling_std))
140
141      goods = (compounds > 0)
142      bads = (compounds < 0)
143      # moods ratio
144      positives = np.count_nonzero(goods)
145      negatives = np.count_nonzero(bads)
146      temp[2].append(positives / len(compounds))
147      temp[3].append(negatives / len(compounds))
148
149      # moods combo
150      consecutive_positives = np.array([sum(1 for _ in group)
151                                        for key, group in groupby(goods) if key == 1])
152      consecutive_negatives = np.array([sum(1 for _ in group)
153                                        for key, group in groupby(bads) if key == 1])
154      temp[4].append((consecutive_positives >= 6).sum())
155      temp[5].append((consecutive_negatives >= 6).sum())
156
157      return temp
158
159  for l_i, l in enumerate(LABEL):
160      print(f"Working on label {l}")
161      reader = pd.read_csv(f"{DIR}/Dataset/MI/{l}.csv",
162                           dtype=dict(zip(columns, [float for i in range(len(columns))])),
163                           chunksize=500000, on_bad_lines='warn', index_col=0,
164                           parse_dates=['timestamp_tweet'])
165
166      if l == 'neg':
167          temp_before_covid = [[] for i in range(len(features))]
```

```python
168             temp_after_covid = [[] for i in range(len(features))]
169
170         temp = [[] for i in range(len(features))]
171         append_df = None
172         i = 0
173         total_user = 0
174         for r in reader:
175             i += 1
176             # Moving chunks
177             if i > 1: r = pd.concat([append_df, r], ignore_index=True, axis=0)
178
179             # Get dataframe without last user
180             users = r['user_id'].unique()
181             last_user = users[-1]
182             mask = (r.user_id == last_user)
183             append_df = r[mask]
184             r = r[~mask]
185             users = users[:-1]
186
187             # Convert timestamp
188             r['date'] = r['timestamp_tweet'].dt.date
189             r['hour'] = r['timestamp_tweet'].dt.hour
190             r.drop(['timestamp_tweet'], axis=1, inplace=True)
191
192             # Group by hours
193             result_hour = r[['tweet_id', 'hour']].groupby(['hour']).count()
194             if len(result_hour) == 0: continue
195             hourly_rate[l_i,:] = result_hour.values.ravel()
196             total_user += len(users)
197
198             # Group by date
199             aggs = dict(zip(columns,['mean' for i in range(len(columns))]))
200             result_date = r.groupby(['user_id', 'date']).agg(aggs)
201             df_before_covid = \
202                     result_date[result_date.index.get_level_values('date') < covid_date]
203             df_after_covid =  \
204                     result_date[result_date.index.get_level_values('date') >= covid_date]
205
206             # Extract features
207             for u in users:
208                 compounds = result_date.loc[u, 'compound']
209                 temp = extract_feature(compounds, temp)
210
211                 # Extract data for covid analysis
212                 if l == 'neg' and
213                     u in df_before_covid.index.get_level_values('user_id') and
214                     u in df_after_covid.index.get_level_values('user_id'):
215                     temp_before_covid = extract_feature(df_before_covid.loc[u, 'compound'],
216                                                         temp_before_covid)
217                     temp_after_covid = extract_feature(df_after_covid.loc[u, 'compound'],
218                                                        temp_after_covid)
219
220             print("Chunk number", i, total_user)
221
222         hourly_rate[l_i,:] /= total_user
223         df = pd.DataFrame([list(i) for i in zip(*temp)], columns=features)
224         dfs.append(df)
225
226         print("FINISHED")
227
228 for i, df in enumerate(dfs): df.to_csv(f"{DIR}/Dataset/MI/{LABEL[i]}_extracted.csv")
229
230 """# Analysis
231
232 ## Window-based feature, moods ratio, moods combo
233 """
234
235 summ = pd.DataFrame(index=features, columns=DISORDERS)
```

```python
236
237 for col in features:
238     for i in range(len(DISORDERS)):
239         t_statistic1, p_value1 = stats.ttest_ind(dfs[i][col].values,
240                                                    dfs[-1][col].values)
241         value = "***" if p_value1 < 0.001 else "**"
242                         if p_value1 < 0.01 else "*"
243                         if p_value1 <0.05 else "."
244                         if p_value1 < 0.1 else "-"
245         color = "r" if round(t_statistic1,10) < 0 else "g"
246                     if round(t_statistic1,10) > 0 else ""
247         summ.at[col, DISORDERS[i]] = value + color
248
249 summ
250
251 # Covid analysis
252 df_before_covid = pd.DataFrame([list(i) for i in zip(*temp_before_covid)], columns=features)
253 df_after_covid = pd.DataFrame([list(i) for i in zip(*temp_after_covid)], columns=features)
254
255 for col in features:
256     t_statistic1, p_value1 = stats.ttest_ind(df_after_covid[col].values,
257                                                df_before_covid[col].values)
258     value = "***" if p_value1 < 0.001 else "**"
259                     if p_value1 < 0.01 else "*"
260                     if p_value1 <0.05 else "."
261                     if p_value1 < 0.1 else "-"
262     color = "r" if round(t_statistic1,10) < 0 else "g"
263                 if round(t_statistic1,10) > 0 else ""
264     summ.loc[col, 'after'] = value + color
265
266 summ
267
268 summ.to_latex()
269
270 """## Hourly rate"""
271
272 for i in range(hourly_rate.shape[0]):
273     plt.plot(hourly_rate[i], label=f'{LABEL[i]}')
274
275 # Set the plot title and labels
276 plt.title('Hourly trends posting in each disorder')
277 plt.ylim(0, 70)
278 plt.xticks(range(24))
279 plt.xlabel('Hour')
280 plt.ylabel('Posting frequency')
281
282 # Show the legend
283 plt.legend()
284
285 # Show the plot
286 plt.show()
```

## Code: Extracting psycholinguistics analysis

```python
1 # -*- coding: utf-8 -*-
2 """LIWC.ipynb
3
4 Automatically generated by Colaboratory.
5
6 Original file is located at
7     https://colab.research.google.com/drive/1aqAlaWpS-GNmXBszJ6acoqFSpAybdb5O
8 """
9
10 !pip install liwc
```

```python
11
12 from google.colab import drive
13 drive.mount('/content/drive')
14
15 DIR = "/content/drive/MyDrive/College (Master)/Semester 3/ \
16         COMP90090 - Natural Language Processing for Health/Assignment/Assignment 3/"
17
18 import os
19 import nltk
20 import liwc
21 import pandas as pd
22 import scipy.stats as stats
23 import statsmodels.api as sm
24 from collections import Counter
25 from nltk.tokenize import word_tokenize
26
27 nltk.download('punkt')
28
29
30 import warnings
31 warnings.filterwarnings("ignore")
32
33 liwcPath = DIR + 'Code/LIWC2015_English_Flat.dic'
34 parse, category_names = liwc.load_token_parser(liwcPath)
35
36 category_used = ['function','ppron', 'i', 'we', 'you', 'shehe', 'they',
37                  'ipron', 'adverb', 'conj', 'adj', 'compare',
38                  'affect', 'posemo', 'negemo', 'anx', 'anger', 'sad',
39                  'social', 'family', 'friend', 'female', 'male',
40                  'cogproc', 'insight', 'cause', 'discrep', 'tentat', 'certain', 'differ',
41                  'percept', 'see', 'hear', 'feel',
42                  'bio', 'body', 'health', 'sexual', 'ingest',
43                  'drives', 'affiliation', 'achieve', 'power', 'reward', 'risk',
44                  'focuspast', 'focuspresent', 'focusfuture',
45                  'work', 'leisure', 'home', 'money', 'relig', 'death']
46
47 DISORDERS = ['ppd', 'mdd', 'ocd', 'ptsd', 'adhd',
48              'bipolar', 'anxiety', 'depression']
49 LABEL = DISORDERS + ['neg']
50
51 """# Feature Extraction"""
52
53 # Store data about user_id, tweet_id, date, sentiment, disorder, covid
54 dfs = []
55
56 for l in ['adhd','neg']:
57     print(f"Working on label {l}")
58     reader = pd.read_csv(f"{DIR}/Dataset/tweets/tweets_{l}.csv", dtype=str,
59                          chunksize=500000, on_bad_lines='warn', lineterminator=";")
60
61     temp = [[] for i in range(len(category_used))]
62
63     i = 0
64     for r in reader:
65         i += 1
66         # if i < 101: continue
67         # if i > 100: break
68         df = r[['tweet_id', 'user_id']]
69
70         counters = []
71         r['text'] = r['text'].fillna("").str.lower()
72         for j, text in enumerate(r['text']):
73             try:
74                 tokens = word_tokenize(text)
75                 counter = Counter(category for token in tokens for category in parse(token))
76                 counters.append(counter)
77             except:
78                 display(r.iloc[j,])
```

```python
 79
 80          counter_df = pd.DataFrame.from_dict(counters).fillna(0).astype(int)
 81          counter_df = counter_df.reindex(columns=category_used, fill_value=0)
 82          final_df = pd.concat([df, counter_df], axis=1)
 83
 84          last_user = final_df.user_id.unique()[-1]
 85          mask = (final_df.user_id == last_user)
 86          append_df = final_df[mask]
 87          final_df = final_df[~mask]
 88
 89          # Aggregate values
 90          result = final_df.groupby(['user_id']).agg(dict(zip(category_used,['mean' for i in range(len(catego
 91          for j in range(len(category_used)):
 92              temp[j] += result[category_used[j]].tolist()
 93
 94          print("Chunk number", i)
 95
 96     df = pd.DataFrame([list(i) for i in zip(*temp)], columns=category_used)
 97     dfs.append(df)
 98     print(" FINISHED")
 99
100 """# Analysis"""
101
102 columns = category_used
103 dfs = []
104
105 for l in LABEL:
106     print(f"Working on label {l}")
107     types = dict(zip(columns,['Int64' for i in range(len(columns))]))
108     types['user_id'] = str
109     reader = pd.read_csv(f"{DIR}/Dataset/LIWC/{l}.csv", dtype=types,
110                          chunksize=500000, on_bad_lines='warn', index_col=0)
111
112     temp = [[] for i in range(len(columns))]
113
114     i = 0
115     append_df = None
116     for r in reader:
117         i += 1
118         # Moving chunks
119         if i > 1: r = pd.concat([append_df, r], ignore_index=True, axis=0)
120         # if i > 10: break
121
122         # Get dataframe without last user
123         r.dropna(subset=['user_id'], inplace=True)
124         r.fillna(0, inplace=True)
125         last_user = r.user_id.unique()[-1]
126         mask = (r.user_id == last_user)
127         append_df = r[mask]
128         r = r[~mask]
129
130         # Aggregate values
131         result = r.groupby(['user_id'])
132                     .agg(dict(zip(columns,['mean' for i in range(len(columns))])))
133         for j in range(len(columns)):
134             temp[j] += result[columns[j]].tolist()
135
136         print("Chunk number", i)
137
138     df = pd.DataFrame([list(i) for i in zip(*temp)], columns=columns)
139     dfs.append(df)
140     print("FINISHED")
141
142 for i, df in enumerate(dfs): df.to_csv(f"{DIR}/Dataset/LIWC/{LABEL[i]}_extracted.csv")
143
144 features = category_used
145 summ = pd.DataFrame(index=features, columns=DISORDERS)
146
```

```
147  for col in features:
148      for i in range(len(DISORDERS)):
149          t_statistic1, p_value1 = stats.ttest_ind(dfs[i][col].values,
150                                                    dfs[-1][col].values)
151          value = "***" if p_value1 < 0.001 else "**"
152                        if p_value1 < 0.01 else "*"
153                        if p_value1 <0.05 else "."
154                        if p_value1 < 0.1 else "-"
155          color = "\cellcolor{red!25}"
156                        if round(t_statistic1,10) < 0 else "\cellcolor{green!25}"
157                        if round(t_statistic1,10) > 0 else ""
158          summ.at[col, DISORDERS[i]] = color + value
159
160  summ
161
162  summ.to_latex()
```

## Code: Model classificatiion

```
1   # -*- coding: utf-8 -*-
2   """Prediction Model Full.ipynb
3
4   Automatically generated by Colaboratory.
5
6   Original file is located at
7       https://colab.research.google.com/drive/1XzNifgy3T_go-5QPmG4fqsRZZqlIC2eZ
8   """
9
10  DIR = "/content/drive/MyDrive/College (Master)/Semester 3/ \
11          COMP90090 - Natural Language Processing for Health/Assignment/Assignment 3/"
12
13  from google.colab import drive
14  drive.mount('/content/drive')
15
16  import nltk
17  import numpy as np
18  import pandas as pd
19  import seaborn as sns
20  import scipy.stats as stats
21  import statsmodels.api as sm
22  import matplotlib.pyplot as plt
23  import matplotlib.dates as mdates
24
25  from sklearn.svm import LinearSVC
26  from sklearn.naive_bayes import MultinomialNB
27  from sklearn.model_selection import GridSearchCV
28  from sklearn.linear_model import LogisticRegression
29  from sklearn.ensemble import RandomForestClassifier
30  from sklearn.model_selection import train_test_split
31  from sklearn.metrics import classification_report, roc_auc_score, confusion_matrix
32
33  import warnings
34  warnings.filterwarnings("ignore")
35
36  """# Read Data"""
37
38  FEATURES = ["MI", "MA", "LIWC"]
39  DISORDERS = ['ppd', 'mdd', 'ocd', 'ptsd', 'adhd',
40               'bipolar', 'anxiety', 'depression']
41  LABEL = DISORDERS + ['neg']
42
43  full = []
44  mlabel, label = [], []
45
```

```python
46  for l in LABEL:
47      ea_disorder = []
48      for f in FEATURES:
49          df = pd.read_csv(f"{DIR}Dataset/{f}/{l}_extracted.csv", index_col=0)
50          ea_disorder.append(df)
51      feat = pd.concat(ea_disorder, axis=1)
52      feat['mlabel'] = l
53      feat['label'] = 0 if l == 'neg' else 1
54      full.append(feat)
55
56  """# Preprocess"""
57
58  data = pd.concat(full, axis=0, ignore_index=True)
59  data.fillna(0, inplace=True)
60  data[data.isna().any(axis=1)]
61
62  train, test = train_test_split(data, stratify=data['mlabel'],
63                                 test_size=0.2, random_state=42)
64
65  X_train = train.drop(['mlabel', 'label'], axis=1)
66  y_train = train['label']
67  my_train = train['mlabel']
68
69  X_test = test.drop(['mlabel', 'label'], axis=1)
70  y_test = test['label']
71  my_test = test['mlabel']
72
73  """# Train"""
74
75  def tune(model, param_grid, X_train, y_train):
76      grid_search = GridSearchCV(model, param_grid, cv=5, scoring='accuracy', n_jobs=-1)
77      grid_search.fit(X_train, y_train)
78      print('* Best hyperparameters:', grid_search.best_params_)
79      print('* Best accuracy score:', grid_search.best_score_)
80
81      return grid_search.best_estimator_
82
83  """## Binary Classification"""
84
85  nb = MultinomialNB()
86  svm = LinearSVC()
87  rf = RandomForestClassifier()
88  lg = LogisticRegression()
89  models = [lg, nb, svm, rf]
90
91  param_grids = [
92      {'C': [0.001, 0.01, 0.1, 1, 10]},
93      {'alpha': [0.001, 0.01, 0.1, 1, 10]},
94      {'C': [0.001, 0.01, 0.1, 1, 10]},
95      {'n_estimators': [50, 100, 200], 'max_depth': [None, 10, 20]}
96  ]
97
98  bgrids = []
99  # Perform hyperparameter tuning for each model
100 for i, model in enumerate(models):
101     param_grid = param_grids[i]
102     print(model.__class__.__name__)
103
104     # For common training data
105     bgrids.append(tune(model, param_grid, X_train, y_train))
106
107 """## Multilabel classification"""
108
109 nb = MultinomialNB()
110 svm = LinearSVC()
111 rf = RandomForestClassifier()
112 lg = LogisticRegression()
113 models = [lg, nb, svm, rf]
```

```python
114
115  param_grids = [
116      {'C': [0.001, 0.01, 0.1, 1, 10], 'multi_class': ['multinomial']},
117      {'alpha': [0.001, 0.01, 0.1, 1, 10]},
118      {'C': [0.001, 0.01, 0.1, 1, 10], 'multi_class': ['ovr']},
119      {'n_estimators': [50, 100, 200], 'max_depth': [None, 10, 20]}
120  ]
121
122  mgrids = []
123  # Perform hyperparameter tuning for each model
124  for i, model in enumerate(models):
125      param_grid = param_grids[i]
126      print(model.__class__.__name__)
127
128      # For common training data
129      mgrids.append(tune(model, param_grid, X_train, my_train))
130
131  """# Evaluate"""
132
133  def evaluate(y, y_pred, y_pred_proba=None):
134      report = classification_report(y, y_pred)
135      print('Classification Report:')
136      print(report)
137
138      if y_pred_proba is not None:
139          auroc = roc_auc_score(y, y_pred_proba)
140          print('AUROC:', round(auroc, 4))
141      print()
142
143      return classification_report(y, y_pred, output_dict=True)
144
145  breport = pd.DataFrame(columns=['Acc', 'P', 'R', 'F1'],
146                         index=['Logistic Regression', 'Naive Bayes',
147                                'Support Vector Machine', 'Random Forest'])
148  mreport = pd.DataFrame(columns=['Acc', 'P', 'R', 'F1'],
149                         index=['Logistic Regression', 'Naive Bayes',
150                                'Support Vector Machine', 'Random Forest'])
151
152  """## Binary classification"""
153
154  for i, model in enumerate(bgrids):
155      print("============", model.__class__.__name__, "============")
156      y_pred = model.predict(X_test)
157      if hasattr(model, "predict_proba"):
158          report = evaluate(y_test, y_pred, model.predict_proba(X_test)[:,1])
159      else:
160          report = evaluate(y_test, y_pred)
161
162      breport.iloc[i,:] = [round(report['accuracy'],2),
163                           round(report['macro avg']['precision'],2),
164                           round(report['macro avg']['recall'],2),
165                           round(report['macro avg']['f1-score'],2)]
166
167  breport
168
169  breport.to_latex()
170
171  """## Multilabel classifiction"""
172
173  for i, model in enumerate(mgrids):
174      print("============", model.__class__.__name__, "============")
175      my_pred = model.predict(X_test)
176      report = evaluate(my_test, my_pred)
177
178      mreport.iloc[i,:] = [round(report['accuracy'],2),
179                           round(report['macro avg']['precision'],2),
180                           round(report['macro avg']['recall'],2),
181                           round(report['macro avg']['f1-score'],2)]
```

```
182
183  mreport.to_latex()
184
185  # Show confustion matrix for highest accuracy score
186  classification_matrix = confusion_matrix(my_test, my_pred)
187  ax = sns.heatmap(classification_matrix, annot=True, cmap='Blues', fmt='d')
188
189  ax.set_xlabel('Predicted')
190  ax.set_ylabel('Actual')
191  ax.set_xticklabels(['PPD', 'MDD', 'OCD', 'PTSD', 'ADHD',
192                      'Bipolar', 'Anxiety', 'Depression', 'Control'], rotation=90)
193  ax.set_yticklabels(['PPD', 'MDD', 'OCD', 'PTSD', 'ADHD',
194                      'Bipolar', 'Anxiety', 'Depression', 'Control'], rotation=0)
```