# Distinguishing PTSD from Anxiety Disorders: A Machine Learning Investigation of Linguistic Patterns in Online Mental Health Communities

**Youngmeen Kim**
Department of Applied Linguistics
Georgia State University
ykim140@gsu.edu

## Abstract

Anxiety disorders and post-traumatic stress disorder (PTSD) affect millions of people worldwide. However, the similarities between these mental disorders often make accurate diagnosis challenging. This study aims to build a better detection model and to analyze language patterns and topics in online mental health communities using a combination of Natural Language Processing (NLP), Machine Learning (ML), topic modeling, and linguistic analysis. The classification model achieved an accuracy rate of 87.50% (F-score of 0.87 and 0.88 for PTSD and Anxiety disorder, respectively). The major topics discussed in the PTSD community included sleep, loneliness, and therapeutic options, often describing physical symptoms and negative feelings toward others and oneself (e.g., anger, exhaustion). Individuals in the anxiety disorder community frequently discussed careers, relationships, nausea, medication, diet, breathing, and disease, often expressing their internal negative perceptions (e.g., fear, worry). These findings may help identify and support individuals with these mental health disorders and provide personalized treatment.

## 1 Introduction

There is a high prevalence of anxiety disorders and post-traumatic stress disorder (PTSD) around the world (Kadri et al., 2007; Kessler et al., 2005; Lépine, 2002). PTSD is characterized by a continuous emotional response to trauma or frightening memories (Brewin et al., 2000), while anxiety disorder is characterized by persistent anxiousness and can take many forms, such as Generalized Anxiety Disorder, Social Anxiety Disorder, and Panic Disorder (American Psychiatric Association et al., 2013). Although anxiety disorder and PTSD share several overlapping symptoms, such as hyper-arousal, trembling, and avoidance behaviors, the Diagnostic and Statistical Manual of Mental Disorders (American Psychiatric Association et al., 2013) distin-guishes PTSD from anxiety disorders by categorizing it under 'Trauma- and Stressor-Related Disorders,' emphasizing its unique manifestation in response to traumatic or stressful events.

Many factors prevent people from receiving professional treatment, such as self-denial, ignorance, financial constraints, accessibility, and the stigma associated with mental illness (Alonso et al., 2018). Alternatively, people often turn to the Internet to access information, share opinions, and seek support (Powell et al., 2011). With online mental health platforms continuing to grow, the mental health field needs to be prepared to understand the behaviors of users within these communities. Thus, exploring online mental health communities is more than just an academic endeavor.

Despite growing interest in the analysis of online mental health communities, there seems to be a research gap. While many ML-based studies have developed classification models (Dwyer et al., 2018; Shatte et al., 2019) and topic modeling approaches (Paul and Dredze, 2014) to identify prevalent topics in online mental health communities, linguistic studies that investigate unique textual features (Pennebaker et al., 2003; Al-Mosaiwi and Johnstone, 2018) may benefit from triangulating with ML methodologies to improve validity, reliability, and practicality. This study addresses this gap by employing a triangulated approach that combines NLP, ML approaches, and linguistic analysis. The following research questions (RQs) are formulated to achieve these goals.

**RQ1**: *Can an ML model accurately classify texts from PTSD and anxiety disorders?*

**RQ2**: *What are the major topics within the PTSD and anxiety disorder online communities?*

**RQ3**: *What are the distinctive language patterns in the texts related to PTSD and anxiety disorder?*

## 2 Methods

### 2.1 Data preprocessing

A dataset consisting of posts from two Reddit subreddits, r/PTSD and r/anxiety, was used for this study. It comprised a total of 50,000 posts, with 25,000 from each subreddit, consisting of 7,752,500 tokens. The Python modules NLTK and spaCy were primarily utilized for data preprocessing. This process included removing unnecessary information (e.g., non-words), expanding contracted forms, lemmatizing, and tokenizing.

### 2.2 Feature selection

ML model features included in this study were n-grams (continuous sequences of words), TF-IDF (Term Frequency-Inverse Document Frequency), and BERT (Bidirectional Encoder Representations from Transformers). For the n-gram and TF-IDF features, uni-, bi-, and tri-grams, as well as TF-IDF window sizes of 1, 2, and 3, were utilized to capture both word- and phrase-level linguistic characteristics. Despite its computational expense, BERT was incorporated due to its ability to generate context-sensitive embeddings that capture the contextual meaning of words.

### 2.3 Classification

The SVM (support vector machine) (Cortes and Vapnik, 1995) model was used for text classification, which involves finding a hyperplane in a high-dimensional space that best divides data points into different groups. SVM was chosen because it has been shown to perform better on complex datasets with non-linear boundaries than other ML models, such as logistic regression and decision tree methods. To train and test the ML algorithm, the data were divided into 80% for training and 20% for testing. Hyperparameter tuning was carried out to select the optimal kernel function, regularization parameter C, and gamma parameter as well as the 10-fold cross-validation. The classification results included overall accuracy, precision, recall, F1 score, the coefficients of each feature, and a confusion matrix.

### 2.4 Topic modeling

To identify the major themes within the PTSD and anxiety disorder communities, this study employed BERTopic (Grootendorst, 2022), which utilizes the BERT model to capture the meanings, contexts, and relationships between words in texts. BERTopic leverages an unsupervised clustering algorithm (e.g., HDBSCAN, UMAP) to discern topics without requiring predefined topic numbers from researchers. Additionally, BERTopic generates more human-interpretable topics by utilizing context information from dense embeddings. The keywords produced by BERTopic were further analyzed to conduct a linguistic analysis.

### 2.5 Linguistic analysis

Linguistic analysis was performed using the results from the machine learning classification and the keywords and topics identified through the topic modeling approach. Specifically, the coefficients of n-gram and TF-IDF features were used from the SVM model classification. These coefficients reveal the strength and direction of the relationship between each feature and the target variable, indicating each feature's contribution to the prediction. Additionally, the topics and keywords generated by BERTopic, derived from BERT embeddings, were incorporated into the linguistic analysis to further understand the textual characteristics of the communities studied.

## 3 Conclusion

The SVM classification model with n-grams, TF-IDF, and BERT features yielded high accuracy (87.50%) with F1-scores above 0.87 in classifying PTSD and anxiety disorder texts. Subsequent BERTopic modeling identified prevalent topics in PTSD (e.g., sleep, family, loneliness, flashback) and Anxiety (e.g., career, relationship, nausea, breathing). Linguistic similarities and differences were examined in the writings of individuals with anxiety disorders and PTSD. In both groups of users, first-person singular pronouns were frequently used to describe their mental and physical symptoms. However, users in the PTSD community commonly expressed exhaustion and anger (e.g., *I don't want, I don't feel*), while anxiety sufferers often expressed uncontrollable fear and persistent worries (e.g., *I can't help, I am not able to*). These findings offer insights for creating better PTSD and anxiety disorder detection models and understanding of the texts and individuals with these conditions. Furthermore, it also illuminates why some individuals might opt for online mental health communities instead of face-to-face treatment from mental health professionals.

## Limitations

Several limitations need to be addressed. First, the incorporation of features that capture the sentimental and semantic meaning of words, such as those offered by LIWC, could have enhanced interpretability. Additionally, Reddit's feature allowing users to post anonymously restricted access to the users' demographic information. Gaining further insights into these disorders could be facilitated by analyzing texts from officially diagnosed patients, including demographic information like gender, age, race, and cultural background. Despite these limitations, the findings of this study hold practical implications for practitioners and researchers interested in understanding individuals with PTSD and anxiety disorders.

## Ethics Statement

This research adheres to the ethical standards in computing, with a focus on prioritizing the public good and human well-being. Efforts have been made to achieve high-quality, reliable results, while also respecting privacy and minimizing harm. The primary aim of this work is to enhance public understanding of PTSD and anxiety disorders, thereby making a responsible contribution to the field of mental health. Recognizing the importance of inclusiveness, efforts have been taken to ensure that the findings support diverse communities and maintain the confidentiality of the identities of data source authors.

## References

Mohammed Al-Mosaiwi and Tom Johnstone. 2018. In an absolute state: Elevated use of absolutist words is a marker specific to anxiety, depression, and suicidal ideation. *Clinical Psychological Science*, 6(4):529–542.

Jordi Alonso, Zhaorui Liu, Sara Evans-Lacko, Ekaterina Sadikova, Nancy Sampson, Somnath Chatterji, Jibril Abdulmalik, Sergio Aguilar-Gaxiola, Ali Al-Hamzawi, Laura H Andrade, et al. 2018. Treatment gap for anxiety disorders is global: Results of the world mental health surveys in 21 countries. *Depression and anxiety*, 35(3):195–208.

DSMTF American Psychiatric Association, American Psychiatric Association, et al. 2013. *Diagnostic and statistical manual of mental disorders: DSM-5*, 5 edition. 5. American psychiatric association Washington, DC.

Chris R Brewin, Bernice Andrews, and John D Valentine. 2000. Meta-analysis of risk factors for posttraumatic stress disorder in trauma-exposed adults. *Journal of consulting and clinical psychology*, 68(5):748.

Corinna Cortes and Vladimir Vapnik. 1995. Support-vector networks. *Machine learning*, 20:273–297.

Dominic B Dwyer, Peter Falkai, and Nikolaos Koutsouleris. 2018. Machine learning approaches for clinical psychology and psychiatry. *Annual review of clinical psychology*, 14:91–118.

Maarten Grootendorst. 2022. Bertopic: Neural topic modeling with a class-based tf-idf procedure. *arXiv preprint arXiv:2203.05794*.

Nadia Kadri, Mohamed Agoub, Samir El Gnaoui, Soumia Berrada, and Driss Moussaoui. 2007. Prevalence of anxiety disorders: a population-based epidemiological study in metropolitan area of casablanca, morocco. *Annals of General Psychiatry*, 6:1–6.

Ronald C Kessler, Patricia Berglund, Olga Demler, Robert Jin, Kathleen R Merikangas, and Ellen E Walters. 2005. Lifetime prevalence and age-of-onset distributions of dsm-iv disorders in the national comorbidity survey replication. *Archives of general psychiatry*, 62(6):593–602.

Jean-Pierre Lépine. 2002. The epidemiology of anxiety disorders: prevalence and societal costs. *Journal of Clinical Psychiatry*, 63:4–8.

Michael J Paul and Mark Dredze. 2014. Discovering health topics in social media using topic models. *PloS one*, 9(8):e103408.

James W Pennebaker, Matthias R Mehl, and Kate G Niederhoffer. 2003. Psychological aspects of natural language use: Our words, our selves. *Annual review of psychology*, 54(1):547–577.

John Powell, Nadia Inglis, Jennifer Ronnie, and Shirley Large. 2011. The characteristics and motivations of online health information seekers: cross-sectional survey and qualitative interview study. *Journal of medical Internet research*, 13(1):e20.

Adrian BR Shatte, Delyse M Hutchinson, and Samantha J Teague. 2019. Machine learning in mental health: a scoping review of methods and applications. *Psychological medicine*, 49(9):1426–1448.