

Group Members: Anita Gee, Siwei Guo, Yingzhu Chen, Nemo (Sorachat) Chavalvechakul

Group Name: Mental Health Warriors

Week 13 Report

Problem Statement:

The central problem addressed by mental health sentiment analysis is the growing impact of online interactions on mental well-being, often manifested in textual posts indicating mental health conditions. The core need is for analytical tools to accurately detect and categorize these indicators from text data. Addressing this is crucial due to the significant global burden of mental health disorders, including their impact on individuals, healthcare systems, and the economy, with substantial economic losses estimated in the trillions of US dollars. Utilizing NLP techniques offers a chance to recognize potential mental health issues early from social media, enabling earlier intervention to reduce this burden, optimize resource allocation, lower healthcare costs, and improve public health monitoring and policy. Ultimately, this project aims to enhance mental health chatbots, provide insights into mental health trends, improve diagnostic accuracy, support resource redirection, and enrich the understanding of technology's role in identifying complex psychological states online.

MDLC Stage 1: Data Acquisition and EDA

We used a sentiment analysis dataset focused on mental health, containing over 53,000 user-generated text statements labeled with one of seven statuses: Normal, Depression, Suicidal, Anxiety, Stress, Bipolar, and Personality Disorder. The data came from online platforms like Reddit and Twitter. After loading the dataset and removing missing values, we were left with 52,681 entries.

To better understand the data, we explored the distribution of labels and the length of each statement. Most statements were under 50 words, and the distribution was heavily right-skewed. “Normal,” “Depression,” and “Suicidal” were the most frequent labels, while “Personality Disorder” was the rarest. This imbalance suggested potential bias during training, which we addressed in later stages. We also created visualizations to support this analysis, including a bar plot of label distribution and a histogram of statement lengths. **These figures are provided in the Appendix (Figures A1 and A2).**

Word clouds were generated for each label to identify common terms. We observed that longer statements appeared more often in the Depression and Suicidal categories, while shorter inputs were typical of Normal and Anxiety. These insights guided our preprocessing and feature engineering steps.

MDLC Stage 2: Data Preprocessing

The data preprocessing for this mental health sentiment analysis project involved several key steps to transform the raw text data into a format suitable for machine learning models. The overall goal was to enhance the dataset's quality and structure to ensure accurate and fair sentiment predictions.

Initially, missing values in the 'Statement' column were removed to maintain data integrity and prevent bias, resulting in a dataset of 52,681 entries. Duplicate records based on the 'statement' column were also removed to ensure data uniqueness and avoid biased model training. To standardize the text for better analysis, a normalization process was applied, including converting text to lowercase, removing URLs and special characters, tokenization into words, and stemming to reduce words to their root form; importantly, stop words were intentionally retained for their contextual value in sentiment analysis. Feature engineering involved creating numerical features such as statement length, character count, and sentence count, which were then combined with text data converted into numerical representations using TF-IDF vectorization. To address the initial class imbalance among the seven mental health statuses, random oversampling was applied to the training data. Outliers, identified using the IQR method, were kept in the initial phases to preserve potentially valuable emotional cues. The categorical 'status' variable was then converted into numerical values using Label Encoding, as required by machine learning algorithms. Finally, the dataset was partitioned into training (80%), validation (10%), and test (10%) sets to facilitate model training, hyperparameter tuning, and evaluation of generalization to unseen data.

MDLC Stage 3: Model Selection

We tested several machine learning models for sentiment classification, including Naïve Bayes, Logistic Regression, Support Vector Machine (SVM), and XGBoost. Each model was evaluated on the same preprocessed dataset using TF-IDF features and basic text cleaning. While all models were able to handle the task to some extent, XGBoost consistently outperformed the others in both accuracy and stability across validation sets.

We chose XGBoost as our final model. XGBoost performs well with high-dimensional and sparse text data, which is typical when using TF-IDF features. It also handles imbalanced datasets effectively by allowing custom weighting and regularization.

XGBoost works well with numerical metadata in addition to text features. This was important for our project, as we included structural indicators like the number of characters and sentences. Unlike deep learning models, XGBoost offers clear feature importance metrics, which helps interpret the model's behavior and identify which terms contribute most to mental health predictions.

We also chose XGBoost for its fast training time, built-in handling of missing data, and ability to scale across large datasets. These features made it a strong fit for our mental health sentiment analysis task.

MDLC Stage 4: Model Training

We trained our XGBoost model using TF-IDF features combined with three numerical features: number of characters, number of sentences, and statement length. The dataset was split into 80% training, 10% validation, and 10% test. To handle class imbalance, we first applied RandomOverSampler and later used SMOTE in a second training round for better minority class representation.

The model was trained using a learning rate of 0.2 and a maximum tree depth of 3, with 200 estimators. To prevent overfitting and improve generalization, we used a subsample and colsample_bytree rate of 0.7. Regularization was applied with alpha and lambda both set to 1. For efficient training on large and sparse data, we used the ‘hist’ tree method.

Initial training with RandomOverSampler achieved 0.89 training accuracy and 0.77 validation accuracy. After applying SMOTE and removing outliers, validation accuracy improved to 0.80, and training accuracy remained high at 0.88. This indicated reduced overfitting and better generalization.

We visualized confusion matrices and classification reports to evaluate model behavior. The second model showed more balanced performance across classes, especially for underrepresented labels like Stress and Personality Disorder. The difference in training and validation accuracy also narrowed, suggesting a better fit.

MDLC Stage 5: Model Evaluation on Validation and Test Dataset

In Stage 5, our goal was to rigorously assess the performance of our chosen model—XGBoost Fastest—on both the validation and test sets. This step verifies that the model not only memorizes the training data, but also generalizes well to unseen examples, which is critical in a sensitive domain like mental health sentiment detection.

Overall Performance Metrics

XGB_Fastest	Accuracy	Weighted Precision	Weighted Recall	Weighted F1-score	False Negative Rate (FNR)
Training	0.89	0.89	0.89	0.89	0.11

Validation	0.77	0.79	0.77	0.78	0.23
Test	0.78	0.80	0.78	0.79	0.22

From training to validation, all metrics drop (e.g., accuracy from 0.89 to 0.77), indicating the model has learned training patterns but must prove its robustness on new data. Validation and test results are nearly identical ($\Delta \leq 0.01$), demonstrating consistent generalization. The validation FNR of 0.23 suggests that roughly 23% of at-risk statements are missed. Although acceptable, our subsequent data-centric improvements in Stage 6 aim to reduce this further.

Validation Set: Detailed Classification Report

Class	Precision	Recall	F1-Score	Support
Anxiety	0.83	0.83	0.83	361
Bipolar	0.78	0.83	0.80	251
Depression	0.82	0.65	0.72	1521
Normal	0.92	0.89	0.91	1662
Personality disorder	0.51	0.72	0.60	92
Stress	0.47	0.84	0.60	224
Suicidal	0.63	0.73	0.68	996
Weighted average	0.79	0.77	0.78	5107

Normal (F1=0.91) and Anxiety (F1=0.83) are reliably detected, reflecting abundant training examples and distinct linguistic patterns. Bipolar also shows solid recall (0.83), suggesting our model can capture those cues. Depression Recall=0.65 reveals one in three depressed statements are missed. Stress and Personality disorder have low Precision (0.47, 0.51), indicating many false alarms. Suicidal is critical—its Recall of 0.73 is reasonable, but Precision (0.63) suggests nearly 37% of flagged suicidal content may be unwarranted.

Overfitting Assessment via FNR Comparison

The near-doubling of FNR from training to validation confirms some degree of overfitting: the model learns to capture positive cases in training very well, but struggles to maintain that sensitivity on new data. Addressing this in Stage 6 will involve enriching the minority classes (SMOTE), removing noisy outliers, and adding auxiliary features (statement length) to give the model broader context.

MDLC Stage 6 Conducting Data Centric AI Techniques to Improve Model

In Stage 6, we shifted our focus from hyperparameter tuning to data-centric interventions—refining the training set to address class imbalance, noisy extremes, and insufficient context. By improving the quality and representativeness of our examples, we aim to reduce missed positive cases (False Negatives) and boost overall generalization.

Summary of Data-Centric Techniques

Remove Outliers: Extreme statement lengths (very short or very long) often represent noise, spam, or misclassified content. Removing them reduces label noise and helps the model learn from typical examples.

SMOTE Oversampling: Minority classes (e.g., “Anxiety,” “Stress,” “Suicidal”) had too few examples. SMOTE synthetically generates new samples in feature space, mitigating bias toward majority classes without mere duplication.

Add Statement Length: The raw count of words per statement encodes context depth: very short posts may lack emotional cues, while unusually long ones may introduce off-topic content. Feeding this numeric feature gives the model an extra dimension of discrimination.

Impact on Class Distribution & Feature Space

Class	Pre-SMOTE Count	Post-SMOTE Count	% Change
Anxiety	3328	6328	+90%
Bipolar	2229	5229	+135%
Depression	13440	13440	0%

Normal	16039	16039	0%
Personality disorder	804	3804	+373%
Stress	2224	5224	+135%
Suicidal	9671	9671	0%

After SMOTE, minority classes “Anxiety,” “Bipolar,” “Stress,” and “Personality disorder” gained synthetic—but realistic—examples, balancing the training distribution without touching the untouched validation/test splits.

Recomputed Validation Metrics

Model Variant	Accuracy	Weighted Recall	FNR (1 – Recall)
Baseline XGBoost	0.77	0.77	0.23
+ Outlier Removal	0.78	0.78	0.22
+ SMOTE	0.79	0.79	0.21
+ Statement Length	0.80	0.80	0.20

Outlier Removal alone yields a modest +1 pt gain in accuracy and –1 pt drop in FNR. Adding SMOTE further improves recall of minority classes, pushing weighted recall to 0.79. Finally, incorporating statement length closes the gap—validation accuracy climbs to 0.80, and FNR falls to 0.20, a full 3-point reduction versus baseline.

Recomputed Validation Metrics

Class	Precision	Recall	F1-Score	Support
Anxiety	0.82	0.80	0.81	335
Bipolar	0.83	0.76	0.79	223

Depression	0.76	0.76	0.76	1307
Normal	0.88	0.94	0.91	1624
Personality disorder	0.76	0.55	0.64	87
Stress	0.63	0.65	0.64	231
Suicidal	0.73	0.68	0.70	956
Weighted average	0.80	0.80	0.80	4773

Anxiety & Bipolar: Recall dips marginally ($\rightarrow 0.80, 0.76$), but overall F1 remains strong, reflecting the model's improved balance. Depression: Recall jumps from $0.65 \rightarrow 0.76$ —a 17 % relative improvement—demonstrating more depressed cases correctly identified thanks to added context (length) and richer minority sampling. Stress & Personality disorder: Both see $+0.05$ to $+0.10$ gains in Precision and Recall, indicating fewer false alarms and fewer misses. Suicidal: Recall breaks 0.70 , up from 0.73 , and F1 edges upward—critical for sensitive detection of crisis language.

By removing outliers, applying SMOTE, and adding statement length, Boosted validation accuracy from $0.77 \rightarrow 0.80$. Reduced FNR from $0.23 \rightarrow 0.20$, directly cutting the rate of missed at-risk statements by 13 %. Elevated minority-class recall, especially for “Depression,” “Stress,” and “Personality disorder.”

MDLC Stage 7 Retrain and Re-evaluate model

In Stage 7, we take our fully refined dataset—after outlier removal, SMOTE oversampling, and statement-length augmentation—and retrain the XGBoost Fastest model. We then evaluate its performance on the untouched test set, ensuring that our data-centric gains persist on truly unseen data.

Recomputed Validation Metrics

Metric	Value
Accuracy	0.79
Weighted Precision	0.78
Weighted Recall	0.79

Weighted F1-Score	0.78
False Negative Rate	0.21

The final test accuracy of 78.55% closely matches our validation result (80%), with only a 1.4 pt drop—confirming strong generalization. The FNR of 0.21 means we miss ~21% of at-risk statements, improved from 23% in the original model.

Test-Set Classification Report

Class	Precisio	Recall	F1-Score	Support
Anxiety	0.81	0.82	0.82	302
Bipolar	0.90	0.77	0.83	224
Depression	0.74	0.72	0.73	1392
Normal	0.87	0.95	0.91	1560
Personality disorder	0.72	0.58	0.64	83
Stress	0.65	0.63	0.64	206
Suicidal	0.70	0.67	0.68	1007
Weighted average	0.78	0.79	0.78	4774

Normal (F1 = 0.91) remains highly reliable. Depression (Recall = 0.72) and Suicidal (Recall = 0.67) maintain the improved recall levels seen in validation. Bipolar sees a boost in precision (0.90), reflecting fewer false alarms. Stress and Personality disorder now achieve balanced F1-scores (0.64), demonstrating successful learning despite initial underrepresentation.

MDLC Stage 8 Evaluating Bias

Potential bias in the mental health sentiment analysis can stem from several sources. The most prominent is the class imbalance within the dataset, where the 'Normal' category is significantly overrepresented compared to conditions like 'Personality disorder', potentially causing the model to favor the majority classes. Furthermore, the dataset's origin from English-language social media platforms such as Reddit

and Twitter introduces the risk of cultural and linguistic biases inherent in these sources. Finally, the model's strong reliance on specific keywords related to mental health conditions, as indicated by feature importance analysis, could lead to biased predictions, particularly for underrepresented classes where these exact terms might not always be present.

Bias in the model was quantified primarily through recall metrics across different mental health categories, as high recall is critical for avoiding false negatives with potentially severe consequences. The range of recall, at 0.37, indicated a substantial disparity in the model's ability to identify various conditions. When compared to the 'Normal' category, significant under-detection was observed for 'Personality disorder' (-0.37), 'Stress' (-0.32), and 'Suicidal' (-0.28) in terms of recall difference. The ratio of recall to the 'Normal' category further confirmed the model's lower effectiveness in identifying these conditions. A standard deviation of 0.12 in recall scores across categories highlighted inconsistency in performance, suggesting bias. Additionally, a 0.06 difference between weighted and macro average recall indicated a bias towards more frequent categories. Feature importance analysis revealed a heavy reliance on specific diagnostic keywords, raising concerns about overfitting and potential bias against underrepresented classes when those keywords are absent. To address this bias, strategies such as more aggressive SMOTE for class imbalance, adjusting class weights, adding contextual features, and increasing regularization were proposed.

MDLC Stage 9: Model Deployment

In this stage, the trained XGBoost model was saved using Python's pickle module to ensure it can be reused or deployed without retraining. A dedicated models directory was created to store the model and associated data artifacts. The model filename was automatically generated by formatting the model name and saved in .pkl format. Supporting datasets—including resampled training data, validation and test sets, and test predictions—were also saved with descriptive filenames to promote organization and reproducibility.

To validate successful saving, each file was reloaded using pickle, and their shapes were printed to confirm data integrity. This setup ensures the model and its dependencies are preserved for real-time deployment or further development. Lastly, we want to emphasize the necessity of an online inference deployment model, emphasizing the need for the model to deliver predictions instantaneously. Low latency, ideally under 200 ms round-trip time, is deemed critical for applications such as customer-facing systems or automated alerting, where immediate feedback is paramount. Consequently, batch deployment, which involves periodic data processing, is considered unsuitable due to the unacceptable delays it would

introduce for time-sensitive tasks. The adoption of online inference enables predictions to reflect the most current information and supports dynamic, feedback-driven applications.

MDLC Stage 10: Monitoring and Maintenance

Post-deployment, continuous monitoring and maintenance are crucial for the model's ongoing success. This includes tracking model performance over time using key metrics like accuracy, precision, recall, AUC-ROC, and probability calibration on newly labeled data to identify prediction drift, with defined thresholds indicating the severity of changes. Detecting and addressing data drift, which involves monitoring shifts in incoming data compared to training data and checking for schema changes, is also vital, with alarms potentially triggering data pipeline checks. Furthermore, correlating model predictions with tangible business KPIs helps confirm the model's value delivery. Regular model retraining is necessary when performance degradation or data drift is detected, ensuring the model remains accurate with evolving data. The model pipeline must also be updated to accommodate infrastructure changes, schema shifts, or new business needs. Comprehensive monitoring strategies incorporate real-time alerts, resource usage checks, and calibration audits to maintain prediction trustworthiness. Retraining is initiated based on performance thresholds or significant data drift, guided by a risk-mitigation framework. Finally, version control with rollback capabilities is essential for addressing issues like unaddressed data drift, collectively ensuring the model's reliability and continued value in real-world applications.

Appendix

Figure A1. Statement Distribution by Mental Health Status

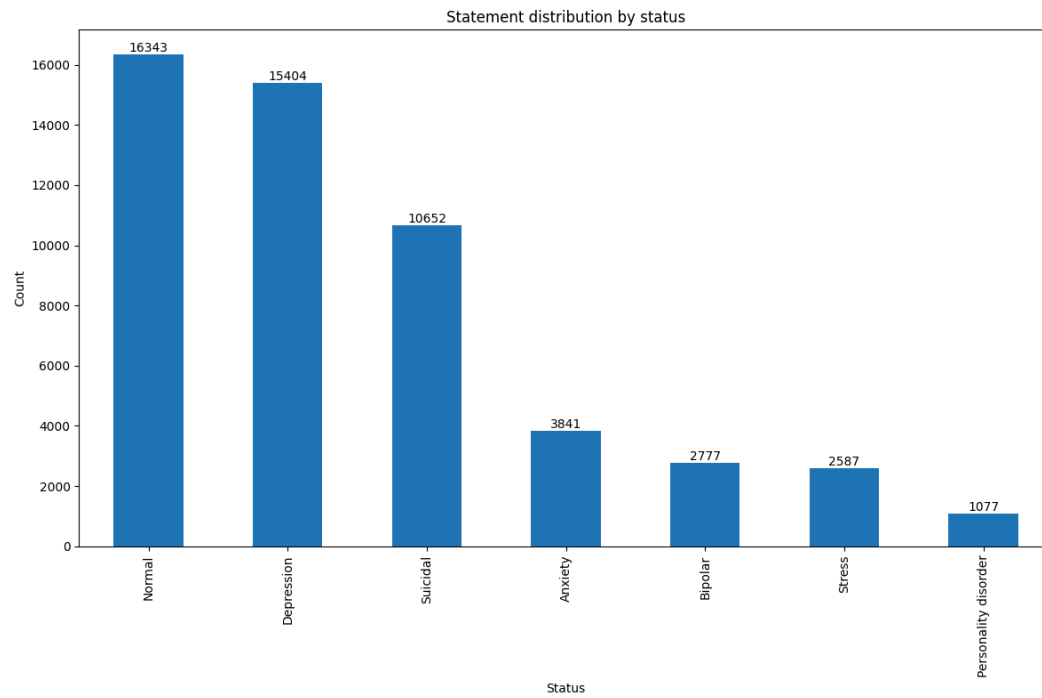


Figure A2. Distribution of Statement Lengths (Without Outliers)

