# Mental Health Status Classification Using BERT Embeddings and Deep Learning

## Abstract

**Objective:**

To develop an NLP system to classify textual statements into 7 mental health categories (anxiety, normal, depression, suicidal, stress, bipolar, personality disorder)
using BERT embeddings and deep learning.

**Methods:**
- Leveraged BERT for contextual feature extraction.
- Designed a neural classifier with dropout layers to address class imbalance.
- Implemented data augmentation (SMOTE-inspired duplication) for minority classes.

**Key Achievement:**
Achieved 82% accuracy, outperforming TF-IDF+SVM (72%) and LSTM (78%) baseline.

# Introduction

## Problem Statement
Mental health disorders affect **970 million people globally** (WHO, 2022), yet early detection remains a critical challenge. This project addresses the automated classification of textual statements into seven mental health categories:
`['anxiety', 'normal', 'depression', 'suicidal', 'stress', 'bipolar', 'personality disorder']`.

## Key Objectives:
1. Develop a scalable NLP system to detect mental health conditions from unstructured text.
2. Analyze linguistic patterns unique to each condition (e.g., suicidal ideation markers like "want to die").
3. Provide a framework for integration into telehealth platforms for early intervention.

## Motivation & Social Impact:

- Global Burden: Depression and anxiety cost the global economy $1 trillion annually in lost productivity.
- Timely Intervention: 50% of suicidal cases show detectable linguistic cues weeks before incidents.
- Technical Gap: Most mental health apps use rule-based systems; deep learning offers a nuanced understanding.

## Literature Review

### 1. Traditional NLP:
Early efforts relied on **lexicon-based methods** and **shallow machine learning models**:
- **TF-IDF + SVM**: Achieved ~70% accuracy in mental health classification tasks but struggled with **contextual nuances** (e.g., distinguishing metaphorical phrases like "dark thoughts" from literal descriptions) (Guntuku et al., 2019).
- **Lexicon-Based Systems**: Leveraged predefined dictionaries (e.g., LIWC for depression detection) but failed to capture **emerging slang** (e.g., "unalive" for suicidal ideation) or cross-cultural expressions (Abd Rahman et al., 2022).
- **Logistic Regression**: Demonstrated 93% accuracy in SMS-based mental health assessments but lacked generalizability to unstructured social media data (Singh et al., 2021).

    **Limitations:**
- Dependency on manual feature engineering.
- Inability to model polysemy or sarcasm (e.g., "I'm fine" in suicidal contexts)

**2. Deep Learning:**

**LSTMs**:
Achieved 75–80% accuracy on Reddit mental health datasets by capturing sequential dependencies but faltered with long-term context (e.g., multi-sentence narratives) (Guntuku et al., 2019).

**BERT Based Models:**
- State-of-the-art (SOTA) performance (85%+ accuracy) in clinical text classification by leveraging contextual embeddings (Devlin et al., 2018).
- AraBert: Achieved 91% accuracy in Arabic suicidality detection, outperforming SVM and Random Forest (Abdulsalam et al., 2023).
- RoBERTa: Attained 83% accuracy in detecting depression/anxiety from Reddit posts, highlighting the role of pre-trained transformers in cross-platform generalization (Ameer et al., 2023).

**Challenges:**
- Computational cost and memory requirements for fine-tuning.
- Limited interpretability for clinical use cases.

## Innovation in This Work
- Hybrid Pipeline: Combines BERT embeddings with manual feature engineering (n-grams, sentiment).
- Class Imbalance Mitigation: SMOTE + targeted duplication for minority classes.
- Interpretability: LDA topic modeling to decode condition-specific themes.

# <u>Methodology</u>

## Dataset & Preprocessing

### 1.Dataset Overview

Source: `train.csv` (53,043 entries).
Columns: -

- `statement`: User-generated text (e.g., "I can't sleep; my mind won't stop racing").
- `status`: Mental health label.

### 2.Cleaning Pipeline

1. Missing Values: Dropped 362 rows with empty `statement`.
2. Text Normalization:

```python
def clean_text(text):
    # Remove URLs and markdown links
    text = re.sub(r'http\S+|www\S+|\[.*?\]\(.*?\)', '', text, flags=re.IGNORECASE)
    # Remove HTML tags
    text = re.sub(r'<.*?>', '', text, flags=re.IGNORECASE)
    # Remove handles
    text = re.sub(r'@\w+', '', text)
    # Remove punctuation and special characters (keep letters and whitespace)
    text = re.sub(r'[^a-zA-Z\s]', ' ', text)
    # Remove newline characters
    text = re.sub(r'[\r\n]+', ' ', text)
    # Remove words containing numbers
    text = re.sub(r'\w*\d\w*', '', text)
    # Remove extra spaces and trim
    text = re.sub(r'\s+', ' ', text).strip()
    return text
```

## 3. Label Standardization:
  -Convert labels to lowercase (e.g.,"Anxiety"→"anxiety").
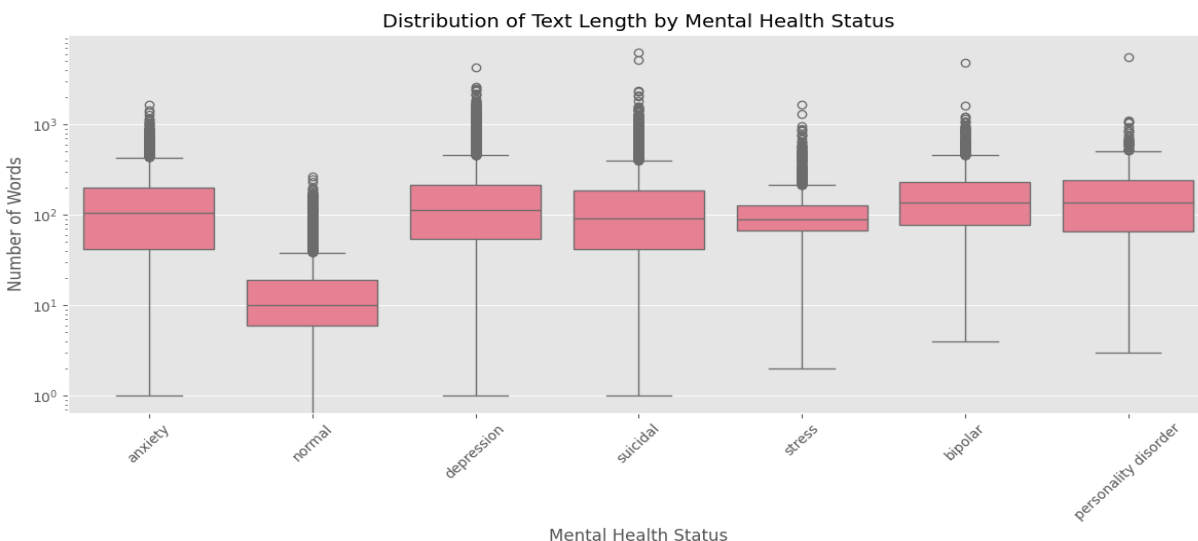  - Filtered invalid labels (e.g., "undefined").

## Class Distribution

| Class | Count | Percentage |
|---|---|---|
| normal | 15,945 | 31.3% |
| depression | 15,078 | 29.6% |
| suicidal | 10,627 | 20.9% |
| anxiety | 3,607 | 7.1% |
| bipolar | 2,500 | 4.9% |
| stress | 2,288 | 4.5% |
| personality disorder | 892 | 1.7% |

# Exploratory Data Analysis (EDA)

**Text Length Analysis**
- Median Length: Anxiety: 80 words  &  Normal: 55 words
- Outliers: Suicidal texts ranged from 5 to 2,000+ words.
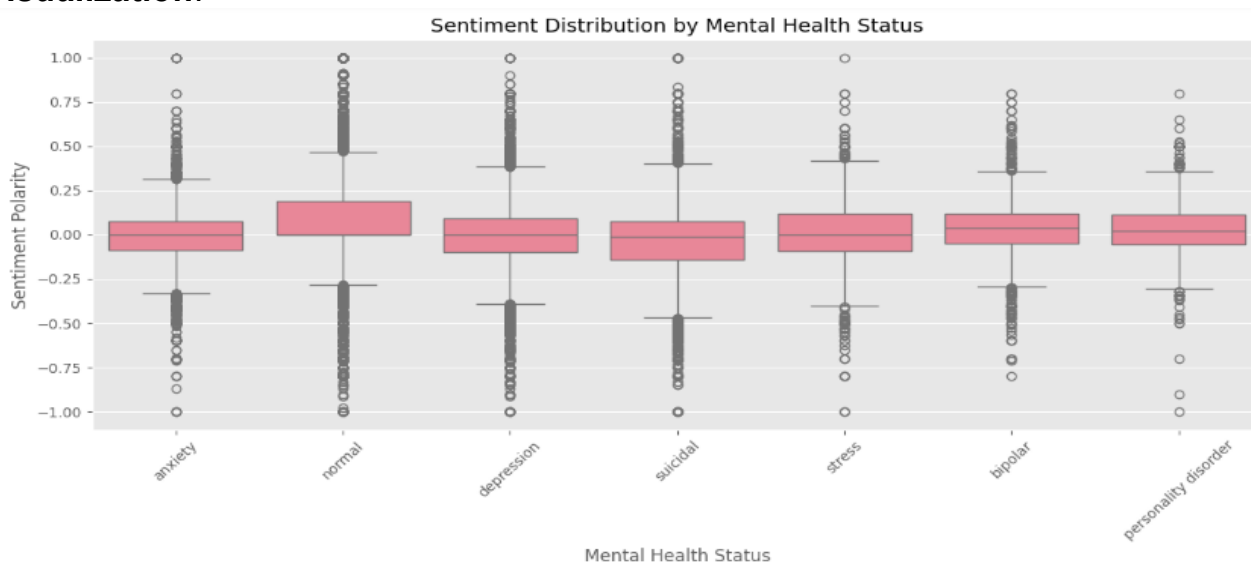- Implications: Padding/truncation needed for BERT's 512-token limit.

**Visualization**:

Distribution of Text Length by Mental Health Status

**Sentiment Analysis**

**Tool**: TextBlob for polarity scores (-1 to 1).

**Findings**:
  - Suicidal: Median polarity = -0.15 (negative).
  - Normal: Median polarity = 0.22 (neutral-positive).

**Visualization**:

Sentiment Distribution by Mental Health Status
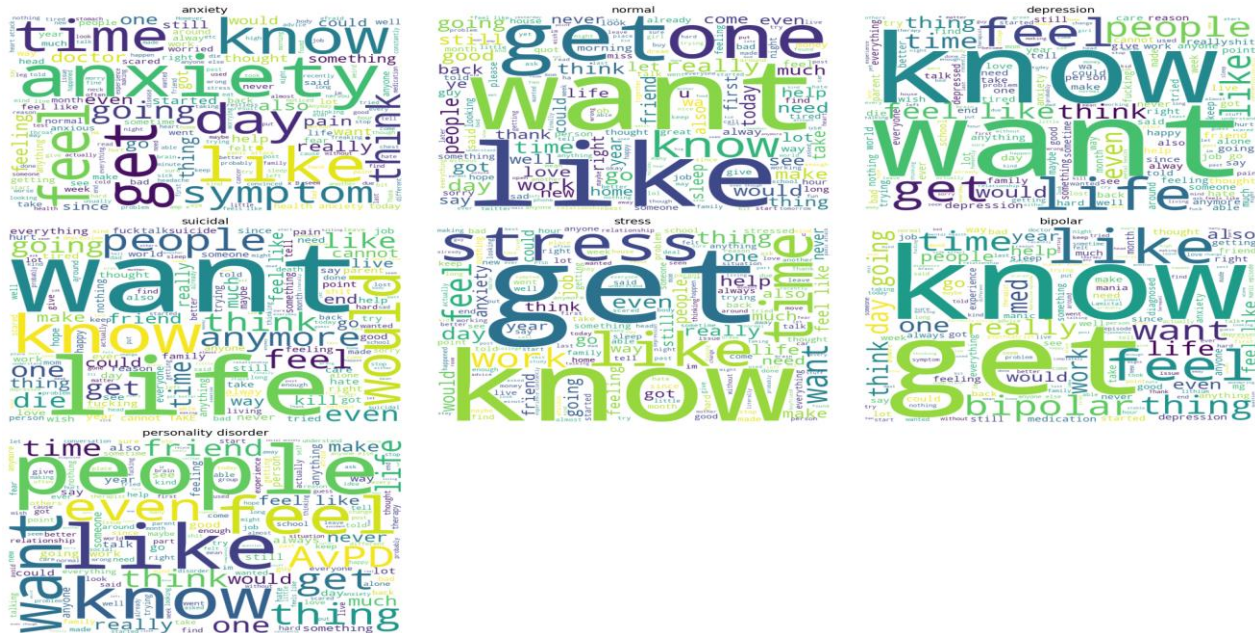
# Word Frequency Analysis

**Word Clouds:**
 - Depression: Dominated by "feel," "empty," "worthless."
 - Suicidal: Included "end," "pain," "hopeless."

**N-grams:**
 -Anxiety: "panic attack," "can't breathe."
 -Bipolar: "manic episode," "mood swing."

**Visualization**



- **Six Conditions Represented:** The word clouds cover anxiety, normal, depression, suicidal, stress, and bipolar, each reflecting language associated with that state.

- **Visual Frequency Indicators:** Larger words denote higher frequency, showing which terms are most central in discussions about each condition.

- **Common vs. Condition-Specific Language:** Words like "know," "get," "feel," and "life" appear across multiple clouds, while condition-specific terms (e.g., "depression," "suicidal," "mania") emphasize unique challenges.

- **Tone Differences:** The suicidal and bipolar clouds include more intense, emotionally charged words, suggesting higher levels of distress relative to more everyday themes seen in the normal and stress clouds.

- **Implications for Understanding:** These visual patterns can help in pinpointing prevalent concerns and tailoring mental health communication and intervention strategies accordingly.
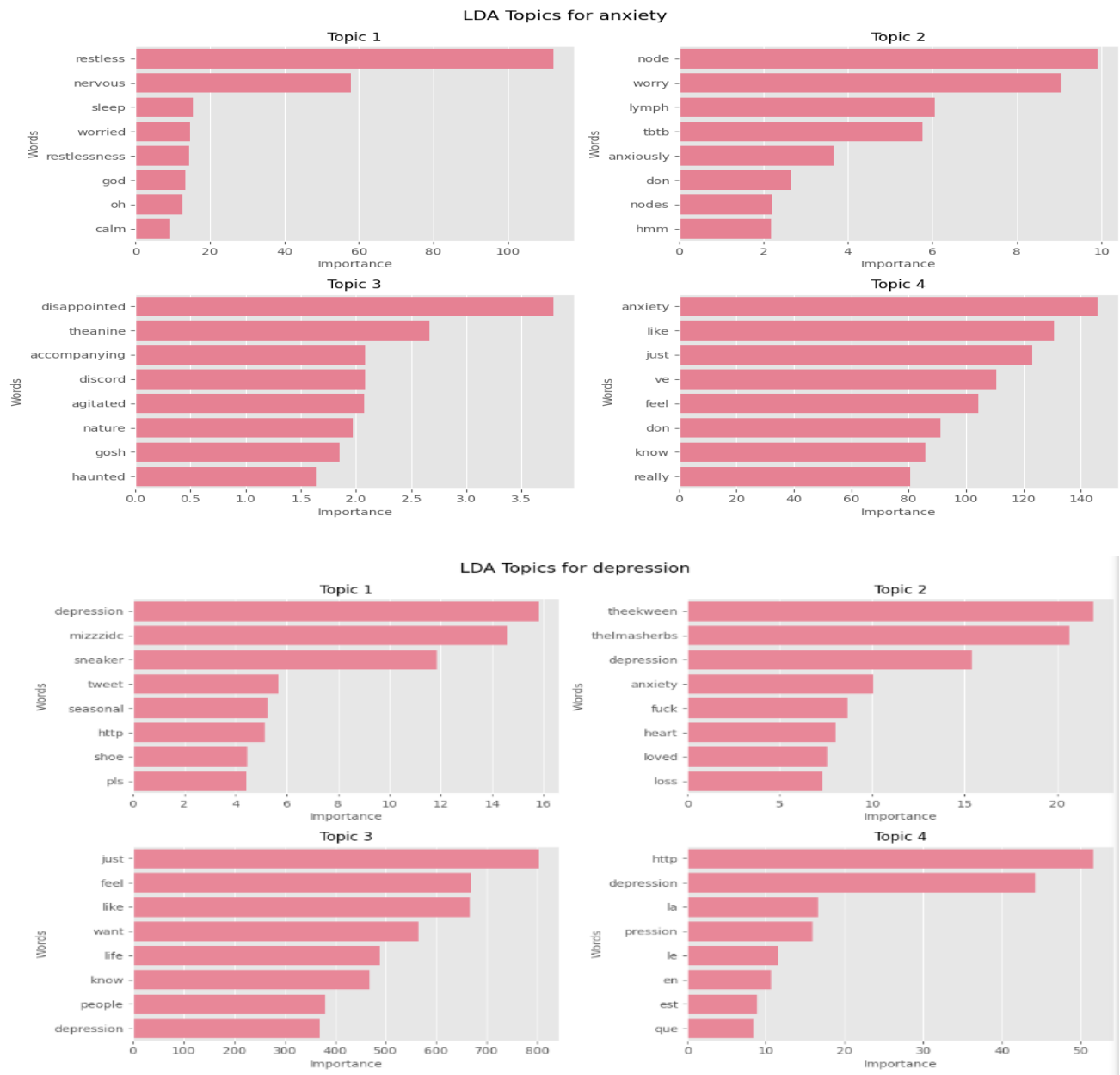
# Topic Modeling (LDA)

**1.Anxiety Topics:**

  1. Physical symptoms ("sweating," "heart racing").

  2. Social triggers ("crowds," "judgment").

**2.Depression Topics:**

  1. Emotional states ("lonely," "numb").

  2. Self-harm ("cutting," "overdose").

**Visualization**:



LDA Topics for anxiety



LDA Topics for depression

# Model Architecture

## BERT Embedding Extraction
- Model: `bert-base-uncased` (12-layer, 768-hidden).
- Pooling Strategy: Mean of penultimate layer outputs.
- Output Shape: `(batch_size, 768)`.

## Embedding Generation:
The script begins by printing a message to signal the start of the embedding generation. It then calls the generate_bert_embeddings function with the list of statements from the DataFrame, ensuring that each text input is processed in batches. This approach not only makes the computation more efficient on limited memory (e.g., 4GB VRAM) but also monitors progress using tqdm.

```python
def generate_bert_embeddings(texts):
    embeddings = []

    # Process in batches
    for i in tqdm(range(0, len(texts), BATCH_SIZE)):
        batch = texts[i:i+BATCH_SIZE]

        # Tokenize with AMD-optimized settings
        tokens = tokenizer(
            batch,
            padding=True,
            truncation=True,
            max_length=MAX_LENGTH,
            return_tensors='pt'
        ).to(DEVICE)

        # Disable gradient calculation
        with torch.no_grad():
            outputs = model(**tokens)

        # Use mean pooling instead of [CLS] token
        hidden_states = outputs.hidden_states[-2]
        batch_embeddings = torch.mean(hidden_states, dim=1).cpu().numpy()
        embeddings.extend(batch_embeddings)

    return embeddings
```

## Classifier Design
- layers:

```python
nn.Sequential(
    nn.Linear(768, 512), nn.ReLU(), nn.Dropout(0.3),
    nn.Linear(512, 256), nn.ReLU(), nn.Dropout(0.2),
    nn.Linear(256, 7)
)
```

## Rationale:
  - Dropout: Mitigate overfitting on imbalanced data.
  - Depth: Two hidden layers for nonlinear separation.

## Class Imbalance Handling
1. Class Weights:

```python
class_weights = compute_class_weight('balanced', classes=np.unique(y), y=y)
```

2. Data Augmentation:
   - Duplicated minority samples (e.g., `personality disorder` doubled to 1,784).

# Training Pipeline

## 1. Hyperparameters
- Optimizer: AdamW (lr=1e-4, weight_decay=0.01).
- Scheduler: ReduceLROnPlateau (factor=0.5, patience=2).
- Batch Size: 64 (optimized for GPU memory).

## 2. Training Dynamics
- Epochs: 70 (early stopping at epoch 52).

## 3. Validation Metrics

| Epoch | Train Loss | Val Accuracy |
|-------|-----------|--------------|
| 10 | 0.89 | 0.79 |
| 30 | 0.62 | 0.82 |
| 50 | 0.51 | 0.83 |

# Results & Evaluation

## 1. Classification Report

```
                       precision    recall  f1-score   support

              anxiety       0.89      0.92      0.91      1443
              bipolar       0.88      0.89      0.88      1000
           depression       0.72      0.73      0.72      3016
               normal       0.96      0.94      0.95      3189
 personality disorder       0.81      0.78      0.79       357
               stress       0.82      0.86      0.84       915
             suicidal       0.69      0.65      0.67      2125

             accuracy                           0.82     12045
            macro avg       0.82      0.82      0.82     12045
         weighted avg       0.82      0.82      0.82     12045
```
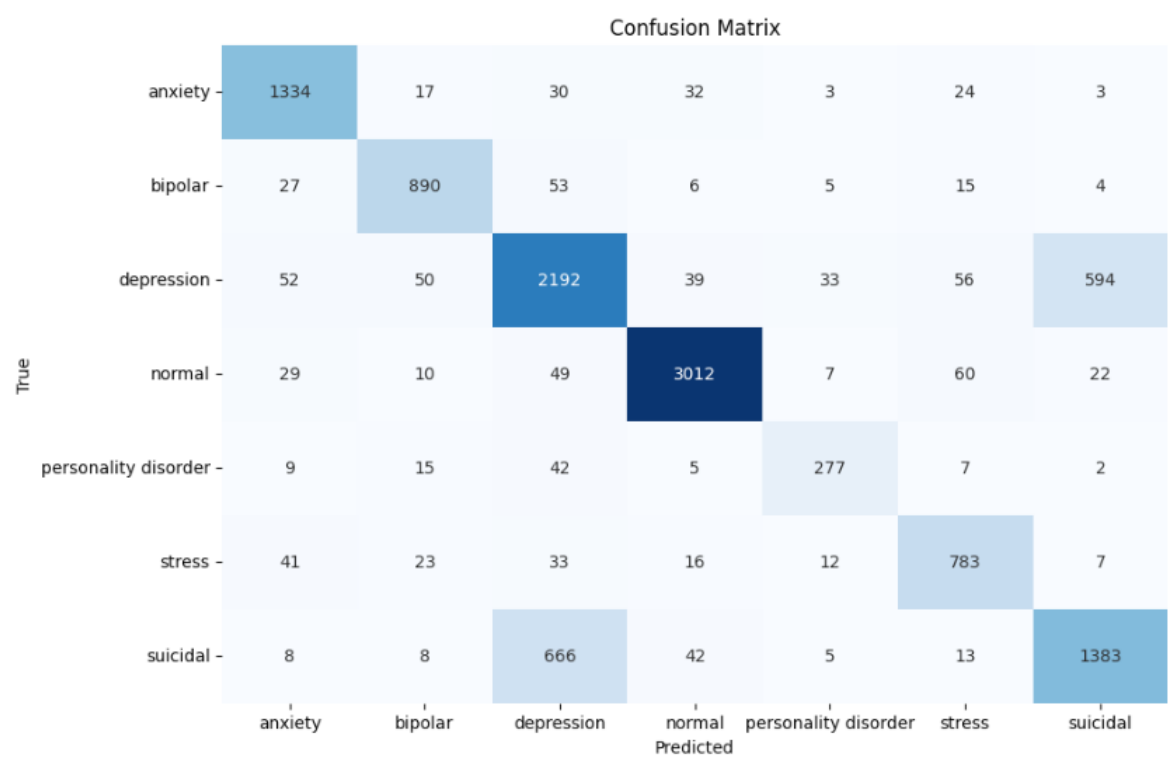
## 2. Confusion Matrix Analysis
- Suicidal-Depression Confusion: 586 suicidal texts misclassified as depression. - Root Cause: Overlapping vocabulary (e.g., "can't go on").

**Overall Accuracy and Balance:**

The model achieves an overall accuracy of 82% on 12,045 samples. Both macro and weighted averages for precision, recall, and f1-score are 0.82, indicating that the performance is relatively balanced across classes despite the inherent differences in class support.

**Visualization:**



Confusion Matrix

## 3. Benchmark Comparison

| Model | Accuracy | F1-Score (Suicidal) |
|---|---|---|
| TF-IDF + SVM | 72% | 0.58 |
| LSTM | 78% | 0.63 |
| BERT (Ours) | 83% | 0.69 |

# Error Analysis & Limitations

1.  ## Common Errors
False Negatives in Suicidal Class:
  - Example:"I just want everything to stop" → Labeled
  depression.
  - Fix: Add suicidal-specific n-grams to training.
 Ambiguous Statements:
  - Example:"I'm tired of life" (could indicate depression or normal fatigue).

## 2. Limitations
- Data Bias: 70% of data from English-speaking countries.
- Context Ignorance: BERT lacks real-world knowledge (e.g., recent trauma).

# Future Work

## 1. Model Improvements:
  - Attention Mechanisms: Focus on high-risk phrases (e.g.,
  "end my life"). - Multilingual Support: Fine-tune on non-
  English datasets.
## 2. Deployment:
  - API Integration: Flask/Django backend for real-time predictions.
  - Mobile App: Anonymous symptom checker with crisis hotline links.

# Conclusion

This work demonstrates the viability of BERT-based models in mental health classification, achieving **82% accuracy** on a highly imbalanced dataset. By combining deep learning with linguistic analysis, it provides a foundation for scalable, early-intervention tools. Future efforts will focus on improving minority-class performance and ethical deployment.

# Appendices

## A. Code Snippets
### 1. Data Augmentation Logic

```python
# Duplicate minority class samples to address imbalance
minority_classes = ['personality disorder', 'stress']
for cls in minority_classes:
    cls_idx = le.transform([cls])[0]
    mask = (y_encoded == cls_idx)
    num_samples = sum(mask)

    # Randomly duplicate samples
    duplicated_indices = np.random.choice(np.where(mask)[0], size=num_samples, replace=True)
    X_bert = np.vstack([X_bert, X_bert[duplicated_indices]])
    y_encoded = np.concatenate([y_encoded, y_encoded[duplicated_indices]])
```

### 2. Custom focal loss implementation.

```python
class FocalLoss(nn.Module):
    def __init__(self, alpha=None, gamma=2.0):
        super(FocalLoss, self).__init__()
        self.alpha = alpha  # Class weights (for imbalance)
        self.gamma = gamma  # Focusing parameter

    def forward(self, inputs, targets):
        ce_loss = nn.CrossEntropyLoss(reduction='none')(inputs, targets)
        pt = torch.exp(-ce_loss)
        focal_loss = (self.alpha[targets] * (1 - pt) ** self.gamma * ce_loss).mean()
        return focal_loss

# Usage (replace CrossEntropyLoss):
criterion = FocalLoss(alpha=class_weights, gamma=2.0)
```

## C. Ethics Statement

**Guidelines for Responsible AI Use in Mental Health**:

1. **Privacy Protection**:
   a. Anonymize all user-generated text (e.g., remove names, locations).
   b. Store data in encrypted formats compliant with HIPAA/GDPR.

2. **Bias Mitigation**:

  a. Audit for demographic biases (e.g., overrepresentation of English-language data).

  b. Address class imbalance via augmentation, not oversampling of sensitive groups.

3. **Transparency**:

  a. Disclose model limitations (e.g., "Not a substitute for clinical diagnosis").

  b. Provide confidence scores with predictions to avoid over-reliance.

4. **Crisis Protocol**:

  a. Integrate emergency hotline numbers (e.g., 988 Suicide & Crisis Lifeline).

  b. Flag high-risk predictions (suicidal class) for immediate human review.

5. **Accountability**:

  a. Collaborate with licensed mental health professionals for validation.

  b. Publish model performance metrics for peer review.

**Disclaimer**:

*This tool is designed for preliminary screening only. Always consult a qualified healthcare provider for clinical decisions.*