

Interim Report – NLP Final Project

# **Analyzing Mental Health: Text Classification for Mental Health Conditions**

Course Lecturer: Dr. Alexander(Sasha) Apartsin

Shahar Saadon | 206560526

Dudi Saadia | 318970944

Shanel Asulin | 205616014

# Project Description

## Data

**Task Type:** Multi-class Text Classification

**Input:** Free-form short English texts (social media statements)

**Output:** One of 7 mental health labels

**Labels:** Normal, Depression, Suicidal, Anxiety, Stress, Bi-Polar, Personality Disorder

**Dataset:** Combined labeled dataset from Kaggle (53,000+ rows)

**Format:** CSV with two columns: statement (text), status (label)

## Challenges

- Informal and varied mental health expressions
- Ambiguous/self-diagnosed language
- Imbalanced classes (e.g., fewer suicidal statements)
- Label noise and missing data

## Evaluation Metrics

- **Accuracy** – Correct predictions rate
- **Precision** – Correctness of positive classifications
- **Recall** – Sensitivity to each class
- **F1-score** – Harmonic mean of precision and recall
- **Confusion Matrix** – Error analysis

# Literature Review

Key Contribution	Best Result	Methods	Dataset & Task	Study
TF-IDF text features, simple yet effective	<b>77%</b> (LightGBM)	Naive Bayes, MLP, LightGBM	10K Reddit posts (6 conditions)	<b>Nova (2023)</b>
Multi-model benchmark with deep/transfer learning	<b>83%</b> (RoBERTa)	Traditional ML, DL, <b>RoBERTa</b>	17K Reddit posts (5 classes)	<b>Ameer et al. (2022)</b>
Detection from <b>general text</b> , not only support forums	<b>81%</b> (RoBERTa)	BERT, XLNet, RoBERTa	100K Reddit posts (9 DSM-5 conditions)	<b>Dinu &amp; Moldovan (2021)</b>

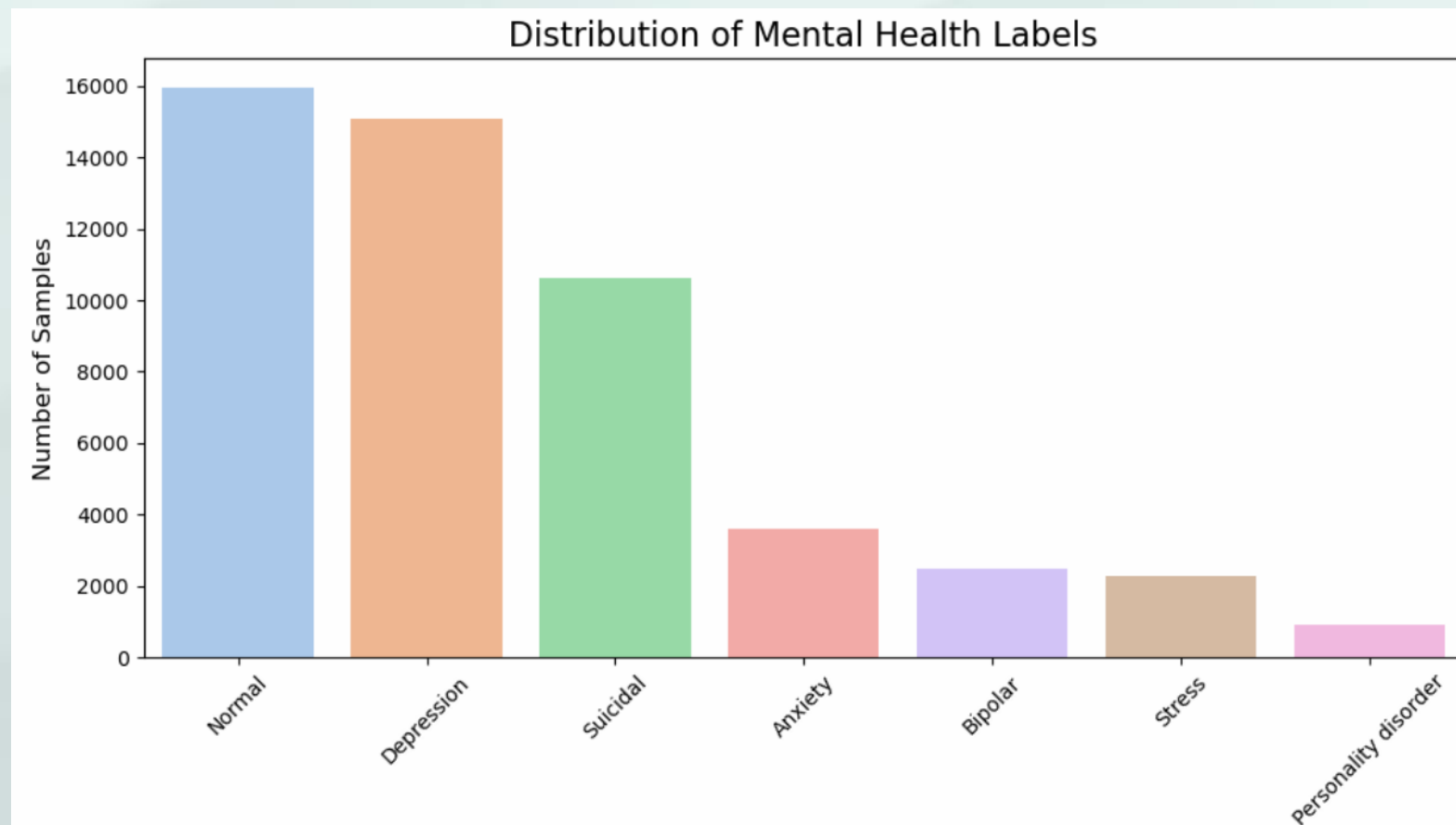
# Work Plan

## Completed so far

- Data preprocessing – loaded, cleaned and verified dataset (53K+ samples).
  - Text vectorization – applied TF-IDF to convert text to numerical features.
  - Baseline model – trained and evaluated Naïve Bayes classifier.
  - Modeling – trained Logistic Regression and Linear SVM models.
  - Evaluation – generated classification reports and confusion matrices.
- 

## Next steps

- Feature optimization – try n-grams, max\_features, etc.
- Advanced modeling – experiment with BERT or RoBERTa.
- Error analysis – inspect misclassified samples, refine labels if needed.
- Visualization – add label distributions, top words, and confusion matrix plots.

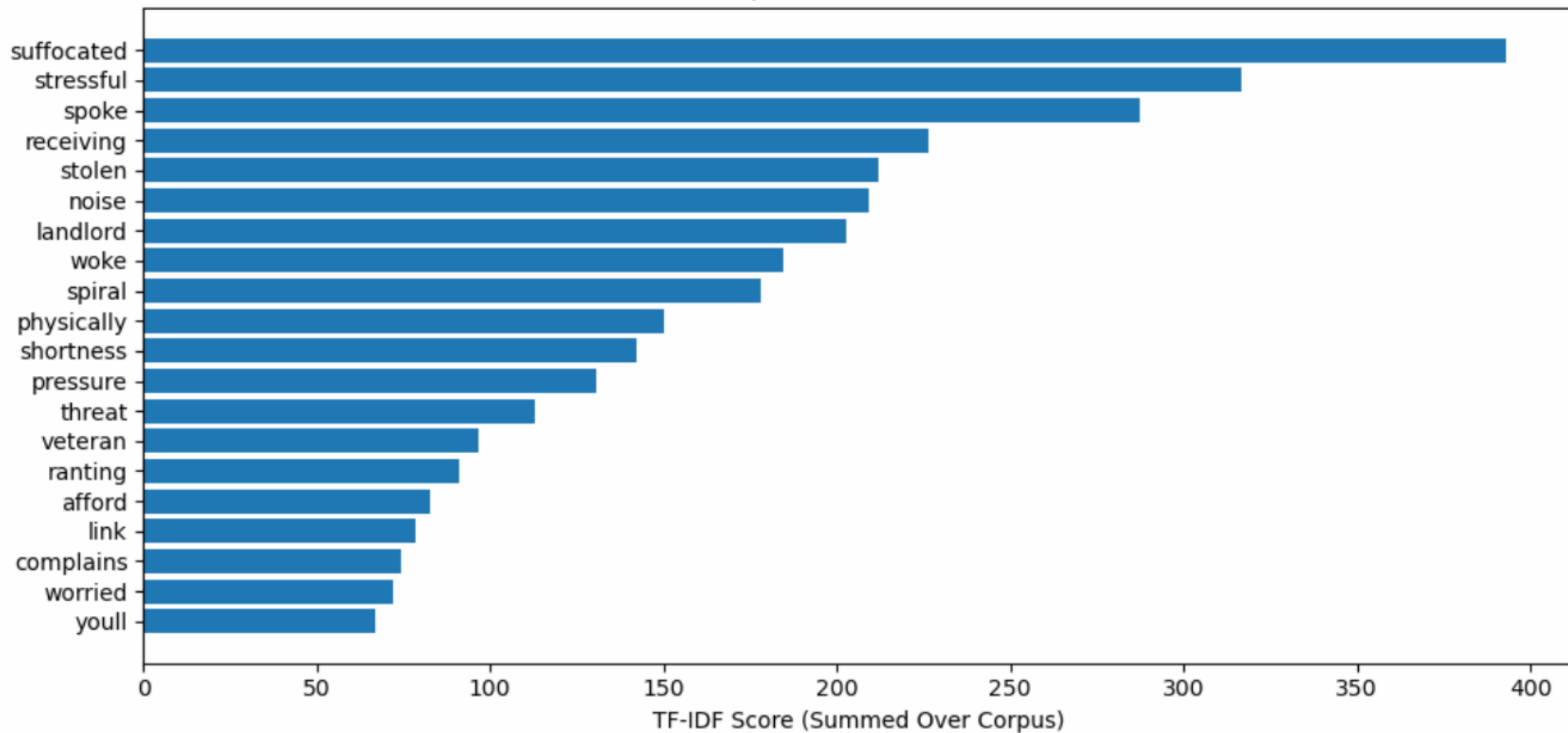


.1

Metric	Naïve bayes	Logistic regression	SVM (LinearSVC)
Accuracy	0.67	0.75	0.74
Marco F1 score	0.51	0.68	0.69
Weighted F1 score	0.65	0.74	0.74

.2

Top 20 TF-IDF Words



# Key Insights & Next Steps

- **Baseline models** like Naive Bayes, Logistic Regression, and SVM performed well on clearer categories (e.g., Normal, Depression) but struggled with more ambiguous ones like Stress and Personality Disorder.
- **Class imbalance** had a noticeable impact on performance, especially in the **Macro F1-score**, emphasizing the importance of using multiple evaluation metrics.
- **TF-IDF helped identify high-impact words**, but it lacks the ability to understand deeper semantic context or word meaning in relation to surrounding words.
- Moving forward, we'll explore **context-aware models like BERT or RoBERTa**, which are better suited for capturing emotional subtleties and complex phrasing.

**Thank You 😊**

---