

Introduction

Group Members: Patrick Xiao, Reiley Vogel, Alexandra Brady

One in five U.S. adults experiences mental illness (cited from Mental Illness – National Institute of Mental Health (NIMH)). Understanding trends in mental health patterns is critical for identifying emerging concerns within specific demographics. Anticipating fluctuations in diagnoses can enable practitioners to provide more effective patient care by allowing time to research new treatment methodologies, explore innovative forms of care, and ensure they are equipped to support individuals with sensitive needs. Accurately prescribing treatments is essential to patient outcomes, but serving as a reliable support system is equally valuable. Early intervention and prevention strategies have the potential to reduce symptom severity or even prevent the onset of mental health conditions altogether. Additionally, healthcare systems can leverage these insights to allocate resources more strategically and efficiently.

Our team utilizes tools such as NLP and sentiment analysis to investigate the potential for predicting mental health diagnoses from clinical notes. This approach aims not only to assist patients in managing their symptoms earlier but also to empower clinicians to anticipate potential cases and ensure that necessary resources are readily available. Our overarching goal is to determine whether artificial intelligence can serve as an efficient and reliable resource for clinicians to better understand the evolving health needs of current and future populations, enabling them to strategically plan and deliver more proactive care.

Overview

Initially, we aimed to analyze anonymized clinical notes, but due to data access limitations, we shifted focus to publicly available data. Our project aims to analyze two key datasets: one focused on mental health treatments, capturing patient details like age, diagnosis, and treatment outcomes, and another exploring sentiment data sourced from multiple platforms including social media posts such as Twitter and Reddit.

For our question, we're asking: **How can we use patient data and sentiment analysis to predict treatment outcomes or identify factors influencing mental health conditions?** Some of the tasks that we will perform in order to answer this includes applying NLP to sentiment data to extract themes using BERT embeddings and topic modeling, building predictive models to correlate patient characteristics and sentiment patterns with treatment outcomes, and conducting some exploratory analysis to identify any critical factors such as stress levels or therapy types.

Why This Matters

Mental health is a growing public health challenge, and traditional diagnostic practices often rely on delayed, reactive approaches. Our project explores the potential of combining patient metadata with NLP-driven sentiment analysis to proactively identify mental health conditions and anticipate treatment outcomes.

The ability to predict diagnoses or treatment outcomes based on text data and patient profiles could revolutionize how we approach mental health care. Our models offer the possibility of earlier intervention, more accurate diagnoses, and optimizes resource allocation-all of which are vital in a system where time and attention are limited. This project helps bridge the gap between unstructured public discourse and clinical insight, potentially enabling healthcare systems to understand and respond to emerging mental health trends in real time.

Datasets Used

Dataset 1 (Sentiment Analysis for Mental Health):

<https://www.kaggle.com/datasets/suchintikasarkar/sentiment-analysis-for-mental-health>

Dataset 2 (Mental Health Diagnosis and Treatment Monitoring):

<https://www.kaggle.com/datasets/uom190346a/mental-health-diagnosis-and-treatment-monitoring>

✓ Data Cleaning

To ensure the quality and consistency of our data, we began by cleaning the sentiment dataset, which contained free-form text statements from social media platforms like Twitter and Reddit. To do this we:

- Removed corrupted entries, such as unnamed columns and rows with missing values
- Applied a custom text-cleaning pipeline that included lowercasing, removing URLs, mentions, hashtags, punctuation, and numbers.
- Removed common stopwords using NLTK and applied lemmatization to reduce each word to its base form (essential to normalize language variation).

```
1 # Libraries
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import seaborn as sns
```

```

6 import re
7 import nltk
8 from nltk.corpus import stopwords
9 from nltk.tokenize import word_tokenize
10 from nltk.stem import WordNetLemmatizer
11 from sklearn.feature_extraction.text import TfidfVectorizer, CountVectorizer
12 from sklearn.decomposition import LatentDirichletAllocation, NMF
13 from sklearn.model_selection import train_test_split
14 from sklearn.metrics import classification_report, confusion_matrix
15 from sklearn.linear_model import LogisticRegression
16 from sklearn.ensemble import RandomForestClassifier
17 from sklearn.svm import SVC
18 from sklearn.preprocessing import LabelEncoder
19 from sklearn.metrics import accuracy_score
20 from wordcloud import WordCloud
21 from textblob import TextBlob
22 import spacy
23 from transformers import BertTokenizer, BertModel
24 import torch
25 from collections import Counter
26 from wordcloud import WordCloud
27 from textblob import TextBlob
28
29 # Download NLTK resources
30 nltk.download('punkt')
31 nltk.download('punkt_tab')
32 nltk.download('stopwords')
33 nltk.download('wordnet')
34 nltk.download('omw-1.4')
35
36 # Configurations
37 nltk.download(['punkt', 'stopwords', 'wordnet'])
38 pd.set_option('display.max_columns', None)
39 sns.set_style('whitegrid')
40 plt.rcParams['figure.figsize'] = (12, 6)
41
42
43 # Set style for visualizations
44 sns.set_style('whitegrid')
45 plt.rcParams['figure.figsize'] = (12, 6)

```

```

[🔄] [nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]   Unzipping tokenizers/punkt.zip.
[nltk_data] Downloading package punkt_tab to /root/nltk_data...
[nltk_data]   Unzipping tokenizers/punkt_tab.zip.
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]   Unzipping corpora/stopwords.zip.
[nltk_data] Downloading package wordnet to /root/nltk_data...
[nltk_data] Downloading package omw-1.4 to /root/nltk_data...
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]   Package punkt is already up-to-date!
[nltk_data] Downloading package stopwords to /root/nltk_data...

```

```
[nltk_data] Package stopwords is already up-to-date!
[nltk_data] Downloading package wordnet to /root/nltk_data...
[nltk_data] Package wordnet is already up-to-date!
```

```
1 # Dataset 1: Mental Health Diagnosis Treatment
2 df_treatment = pd.read_csv('mental_health_diagnosis_treatment_.csv')
3
4 # Dataset 2: Sentiment
5 df_sentiment = pd.read_csv('Combined Data.csv')
```

```
1 df_treatment.head()
```



	Patient ID	Age	Gender	Diagnosis	Symptom Severity (1-10)	Mood Score (1-10)	Sleep Quality (1-10)	Physical Activity (hrs/week)	Medication
0	1	43	Female	Major Depressive Disorder	10	5	8	5	Mc Stabiliz
1	2	40	Female	Major Depressive Disorder	9	5	4	7	Antipsycho
2	3	55	Female	Major Depressive Disorder	6	3	4	3	SS

```
1 df_sentiment.head()
```



	Unnamed: 0	statement	status
0	0	oh my gosh	Anxiety
1	1	trouble sleeping, confused mind, restless hear...	Anxiety
2	2	All wrong, back off dear, forward doubt. Stay ...	Anxiety
3	3	I've shifted my focus to something else but I'...	Anxiety
4	4	I'm restless and restless, it's been a month n...	Anxiety

```
1 df_treatment.isna().sum()
```



	0
Patient ID	0
Age	0
Gender	0
Diagnosis	0
Symptom Severity (1-10)	0
Mood Score (1-10)	0
Sleep Quality (1-10)	0
Physical Activity (hrs/week)	0
Medication	0
Therapy Type	0
Treatment Start Date	0
Treatment Duration (weeks)	0
Stress Level (1-10)	0
Outcome	0
Treatment Progress (1-10)	0
AI-Detected Emotional State	0
Adherence to Treatment (%)	0

dtype: int64

```
1 df_sentiment.isna().sum()
```



	0
Unnamed: 0	0
statement	362
status	0

dtype: int64

```
1 df_sentiment.dropna(inplace=True)
```

```
1 df_sentiment = df_sentiment.drop(columns=['Unnamed: 0'])
```

```
1 df_sentiment.isna().sum()
```

↵

	0
Unnamed: 0	0
statement	0
status	0

dtype: int64

```
1 print(df_treatment.info())
2 print(df_sentiment.info())
```

↵

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 500 entries, 0 to 499
Data columns (total 17 columns):

#	Column	Non-Null Count	Dtype
0	Patient ID	500 non-null	int64
1	Age	500 non-null	int64
2	Gender	500 non-null	object
3	Diagnosis	500 non-null	object
4	Symptom Severity (1-10)	500 non-null	int64
5	Mood Score (1-10)	500 non-null	int64
6	Sleep Quality (1-10)	500 non-null	int64
7	Physical Activity (hrs/week)	500 non-null	int64
8	Medication	500 non-null	object
9	Therapy Type	500 non-null	object
10	Treatment Start Date	500 non-null	object
11	Treatment Duration (weeks)	500 non-null	int64
12	Stress Level (1-10)	500 non-null	int64
13	Outcome	500 non-null	object
14	Treatment Progress (1-10)	500 non-null	int64
15	AI-Detected Emotional State	500 non-null	object
16	Adherence to Treatment (%)	500 non-null	int64

dtypes: int64(10), object(7)

memory usage: 66.5+ KB

None

<class 'pandas.core.frame.DataFrame'>
Index: 52681 entries, 0 to 53042
Data columns (total 3 columns):

#	Column	Non-Null Count	Dtype
0	Unnamed: 0	52681 non-null	int64
1	statement	52681 non-null	object
2	status	52681 non-null	object

dtypes: int64(1), object(2)

memory usage: 1.6+ MB

None

```
1 import re
2 import string
```

```

3 import nltk
4 from nltk.corpus import stopwords
5 from nltk.stem import WordNetLemmatizer
6
7 stop_words = set(stopwords.words('english'))
8 lemmatizer = WordNetLemmatizer()
9
10 def clean_text(text):
11     # Lowercase
12     text = text.lower()
13     # Remove URLs
14     text = re.sub(r"http\S+|www\S+|https\S+", '', text, flags=re.MULTILINE)
15     # Remove mentions and hashtags
16     text = re.sub(r'\@w+|\#', '', text)
17     # Remove punctuation
18     text = text.translate(str.maketrans('', '', string.punctuation))
19     # Remove numbers
20     text = re.sub(r'\d+', '', text)
21     # Remove stopwords and lemmatize
22     words = text.split()
23     words = [lemmatizer.lemmatize(word) for word in words if word not in stop_words]
24     # Join back into a single string
25     return " ".join(words)
26
27 # Create a new column with cleaned text
28 df_sentiment['clean_statement'] = df_sentiment['statement'].apply(clean_text)
29
30 # Preview cleaned text
31 df_sentiment[['statement', 'clean_statement']].sample(5)
32

```

 [Show hidden output](#)

✓ Exploratory Data Analysis

```

1 # Overview of demographics and treatment
2 sns.histplot(df_treatment['Age'], bins=20, kde=True)
3 plt.title("Age Distribution")
4 plt.show()
5
6 sns.countplot(data=df_treatment, x='Gender')
7 plt.title("Gender Distribution")
8 plt.show()
9
10 sns.countplot(data=df_treatment, y='Diagnosis', order=df_treatment['Diagnosis'].
11 plt.title("Diagnosis Frequencies")
12 plt.show()
13

```


```
14 # Treatment outcome by therapy type
15 sns.countplot(data=df_treatment, x='Therapy Type', hue='Outcome')
16 plt.title("Treatment Outcome by Therapy Type")
17 plt.xticks(rotation=45)
18 plt.show()
19
```

[Show hidden output](#)

```
1 # Sentiment class distribution
2 sns.countplot(data=df_sentiment, x='status')
3 plt.title("Sentiment Status Distribution")
4 plt.show()
5
6 # Word cloud by status
7 for status in df_sentiment['status'].unique():
8     text = " ".join(df_sentiment[df_sentiment['status'] == status]['statement'])
9     wordcloud = WordCloud(width=800, height=400, background_color='white').gener
10     plt.figure(figsize=(10,5))
11     plt.imshow(wordcloud, interpolation='bilinear')
12     plt.axis("off")
13     plt.title(f"WordCloud for {status} statements")
14     plt.show()
15
```




Sentiment Status	Count
Anxiety	3800
Normal	16200
Depression	15400
Suicidal status	10700
Stress	2600
Bipolar	2800
Personality disorder	1000

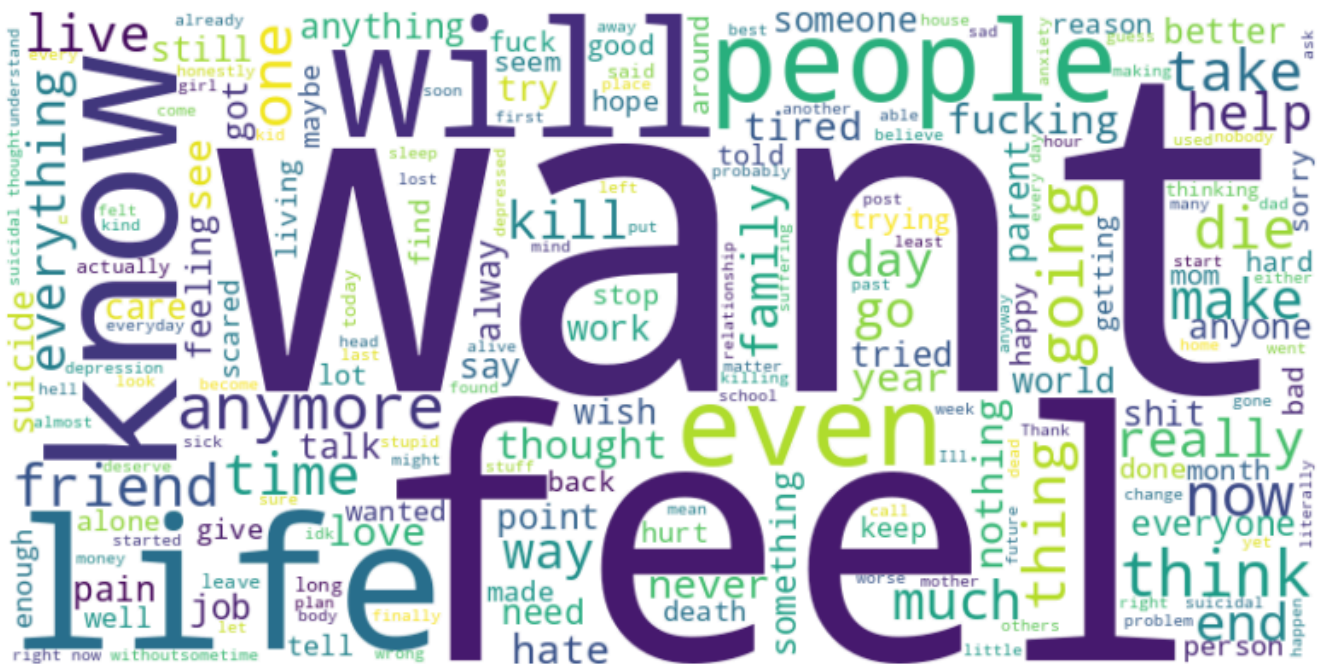
[illegible]



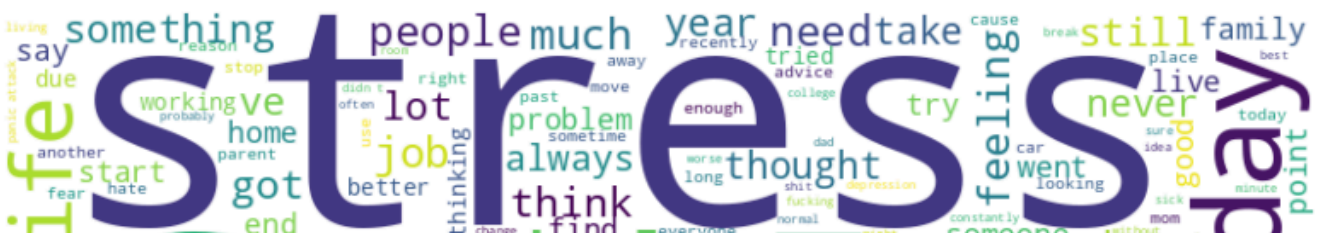
WordCloud for Depression statements

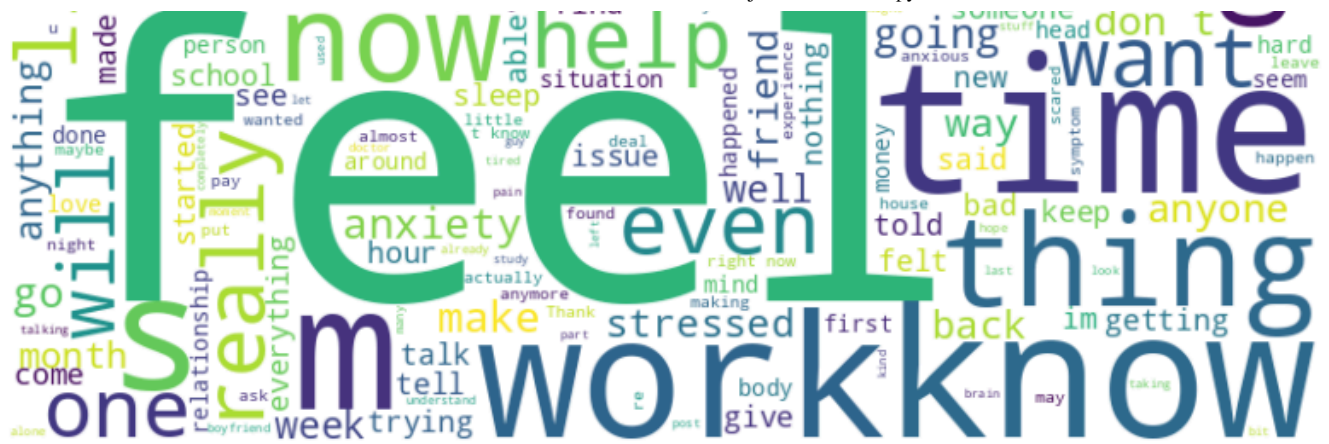


WordCloud for Suicidal statements



WordCloud for Stress statements





WordCloud for Bipolar statements



WordCloud for Personality disorder statements



✓ Named Entity Recognition (NER)

```

1 nlp = spacy.load("en_core_web_sm")
2
3 sample_texts = df_sentiment['statement'].sample(5, random_state=1)
4 for text in sample_texts:
5     doc = nlp(text)
6     print(f"Text: {text}")
7     for ent in doc.ents:
8         print(f"{ent.text} - {ent.label_}")
9     print("-" * 80)
10

```

⇒ Text: i m and i have bad anxiety debilitating i haven t been able to keep a job
early morning hour - TIME
the hard morning - TIME
don - PERSON

Text: Jessica starred in the musical "Legally Blonde" as Elle Woods, the female
Jessica - PRODUCT
Legally Blonde - WORK_OF_ART
Elle Woods - PERSON

Text: I'm so tired I just don't see a point to my suffering, I don't understand

Text: My life 1 year ago was completely different. I was such a chick magnet and
1 year ago - DATE
this past year - DATE
10 years - DATE
one - CARDINAL

Text: RT @no_onespecixl: Know one enjoys my company and I just make everyone mis

✓ Topic Modeling with LDA

```

1 # Vectorize using TF-IDF
2 tfidf = TfidfVectorizer(stop_words='english', max_features=5000)
3 tfidf_matrix = tfidf.fit_transform(df_sentiment['statement'])
4
5 # Fit LDA
6 lda = LatentDirichletAllocation(n_components=5, random_state=42)
7 lda.fit(tfidf_matrix)
8
9 # Display top words per topic
10 def display_topics(model, feature_names, no_top_words):
11     for idx, topic in enumerate(model.components_):
12         print(f"Topic {idx}:")

```

```

13     print(" | ".join([feature_names[i] for i in topic.argsort()[:-no_top_wor
14     print()
15
16 display_topics(lda, tfidf.get_feature_names_out(), 10)
17

```



Topic 0:
just | like | feel | know | life | want | people | time | really | things

Topic 1:
job | work | money | school | yes | just | want | don | like | life

Topic 2:
want | just | life | die | like | feel | tired | fucking | know | anymore

Topic 3:
http | depression | öy | com | rt | morning | don | good | twitter | https

Topic 4:
anxiety | ve | like | just | don | feel | sleep | doctor | day | heart

✓ Sentiment Classification Model

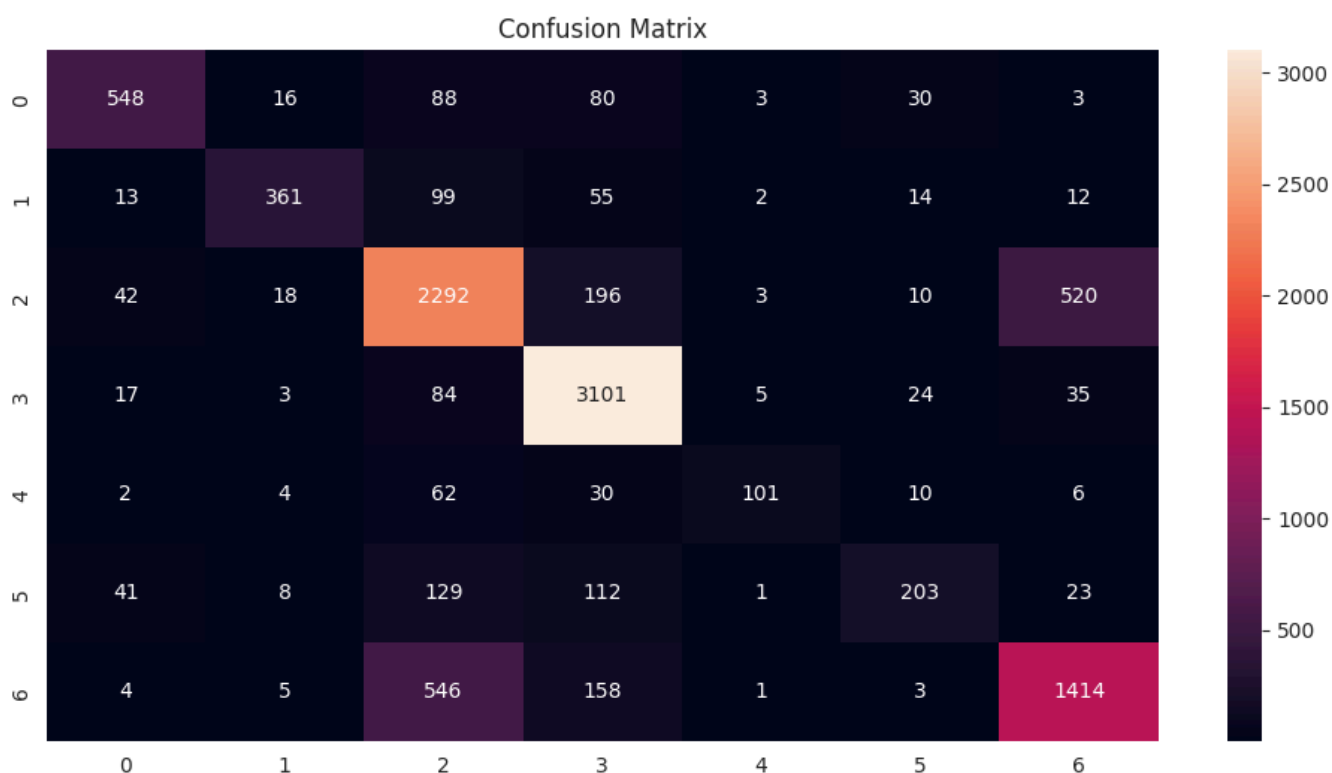
```

1 from sklearn.linear_model import LogisticRegression
2
3 # Prepare data
4 X = tfidf_matrix
5 y = df_sentiment['status']
6
7 # Split
8 X_train, X_test, y_train, y_test = train_test_split(X, y, stratify=y, test_size=
9
10 # Train model
11 clf = LogisticRegression(max_iter=200)
12 clf.fit(X_train, y_train)
13
14 # Evaluate
15 y_pred = clf.predict(X_test)
16 print(classification_report(y_test, y_pred))
17 sns.heatmap(confusion_matrix(y_test, y_pred), annot=True, fmt='d')
18 plt.title("Confusion Matrix")
19 plt.show()
20

```



	precision	recall	f1-score	support
Anxiety	0.82	0.71	0.76	768
Bipolar	0.87	0.65	0.74	556
Depression	0.69	0.74	0.72	3081
Normal	0.83	0.95	0.89	3269
Personality disorder	0.87	0.47	0.61	215
Stress	0.69	0.39	0.50	517
Suicidal	0.70	0.66	0.68	2131
accuracy			0.76	10537
macro avg	0.78	0.65	0.70	10537
weighted avg	0.76	0.76	0.75	10537



▼ Trends

```

1 # Average treatment progress by diagnosis
2 df_progress = df_treatment.groupby("Diagnosis")["Treatment Progress (1-10)"].mean
3 df_progress.plot(kind='barh', title="Average Treatment Progress by Diagnosis", f
4 plt.xlabel("Avg Progress")
5 plt.show()

```

```


6
7 # Correlation heatmap
8 plt.figure(figsize=(10,6))
9 sns.heatmap(df_treatment.select_dtypes(include='number').corr(), annot=True, cma
10 plt.title("Correlation Matrix of Numerical Features")
11 plt.show()
12

```


 [Show hidden output](#)

✓ SpaCy Prediction Model

```
1 print(df_sentiment['status'].unique())
```

 ['Anxiety' 'Normal' 'Depression' 'Suicidal' 'Stress' 'Bipolar'
'Personality disorder']

```
1 print(df_treatment['Diagnosis'].unique())
```

 ['Major Depressive Disorder' 'Panic Disorder' 'Generalized Anxiety'
'Bipolar Disorder']

```

1 status_to_diagnosis = {
2     'Depression': 'Major Depressive Disorder',
3     'Suicidal': 'Major Depressive Disorder',
4     'Anxiety': 'Generalized Anxiety',
5     'Stress': 'Generalized Anxiety',
6     'Bipolar': 'Bipolar Disorder',
7     'Personality disorder': 'Panic Disorder', # Best-guess mapping
8     'Normal': None # We'll exclude these
9 }
10

```

```

1 df_sentiment['Mapped Diagnosis'] = df_sentiment['status'].map(status_to_diagnosi
2 df_sentiment = df_sentiment[df_sentiment['Mapped Diagnosis'].notnull()]

```

```

1 import spacy
2 from sklearn.feature_extraction.text import TfidfVectorizer
3 from sklearn.linear_model import LogisticRegression
4 from sklearn.model_selection import train_test_split
5 from sklearn.metrics import classification_report
6
7 nlp = spacy.load("en_core_web_sm")
8
9 def spacy_tokenizer(text):
10     doc = nlp(text)

```



```

11     return [token.lemma_ for token in doc if not token.is_stop and not token.is_
12
13 vectorizer = TfidfVectorizer(tokenizer=spacy_tokenizer, max_features=5000)
14
15 X = vectorizer.fit_transform(df_sentiment['clean_statement'])
16 y = df_sentiment['Mapped Diagnosis']
17
18 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_
19
20 clf = LogisticRegression(max_iter=1000)
21 clf.fit(X_train, y_train)
22
23 y_pred = clf.predict(X_test)
24 print(classification_report(y_test, y_pred))
25

```

➔ /usr/local/lib/python3.11/dist-packages/sklearn/feature_extraction/text.py:517: warnings.warn(

	precision	recall	f1-score	support
Bipolar Disorder	0.90	0.66	0.77	554
Generalized Anxiety	0.82	0.75	0.78	1256
Major Depressive Disorder	0.91	0.97	0.94	5221
Panic Disorder	0.95	0.48	0.64	237
accuracy			0.89	7268
macro avg	0.89	0.72	0.78	7268
weighted avg	0.89	0.89	0.89	7268

import numpy as np import matplotlib.pyplot as plt import seaborn as sns

✓ Get class labels and their coefficients

```
feature_names = vectorizer.get_feature_names_out() coefs = clf.coef_ labels = clf.classes_
```

Plot top N words for each diagnosis

```

def plot_top_words(coefs, feature_names, labels, top_n=10):
    for idx, label in enumerate(labels):
        top_features = np.argsort(coefs[idx])[-top_n:]
        plt.figure(figsize=(8, 4))
        sns.barplot(x=coefs[idx][top_features], y=[feature_names[i] for i in top_features], palette="viridis")
        plt.title(f"Top {top_n} words for '{label}'")
        plt.xlabel("Coefficient")
        plt.ylabel("Word")
        plt.tight_layout()
        plt.show()

plot_top_words(coefs, feature_names, labels, top_n=10)

```

✓ SpaCy Model with XGBoost Classifier

```

1 from xgboost import XGBClassifier
2 from sklearn.metrics import accuracy_score
3
4 label_encoder = LabelEncoder()
5
6 # Fit the encoder to the 'y_train' labels and transform both 'y_train' and 'y_te
7 y_train_encoded = label_encoder.fit_transform(y_train)
8 y_test_encoded = label_encoder.transform(y_test)
9
10 xgb_clf = XGBClassifier(use_label_encoder=False, eval_metric='mlogloss', random_
11 xgb_clf.fit(X_train, y_train_encoded)
12
13 y_pred_xgb_encoded = xgb_clf.predict(X_test)
14
15 y_pred_xgb = label_encoder.inverse_transform(y_pred_xgb_encoded)
16
17 print("XGBoost Accuracy:", accuracy_score(y_test, y_pred_xgb))
18 print("\nXGBoost Classification Report:\n", classification_report(y_test, y_pred
19

```

➔ /usr/local/lib/python3.11/dist-packages/xgboost/core.py:158: UserWarning: [00:03
Parameters: { "use_label_encoder" } are not used.

```
warnings.warn(msg, UserWarning)
XGBoost Accuracy: 0.908778205833792
```

XGBoost Classification Report:

	precision	recall	f1-score	support
Bipolar Disorder	0.88	0.73	0.80	554
Generalized Anxiety	0.83	0.79	0.81	1256
Major Depressive Disorder	0.93	0.97	0.95	5221
Panic Disorder	0.96	0.59	0.73	237
accuracy			0.91	7268
macro avg	0.90	0.77	0.82	7268
weighted avg	0.91	0.91	0.91	7268

✓ Comparison of SpaCy and XGBoost

```

1 from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay
2 import matplotlib.pyplot as plt
3
4 # Train the Logistic Regression model
5 from sklearn.linear_model import LogisticRegression
6 clf = LogisticRegression(max_iter=1000, random_state=42)

```

```

7 clf.fit(X_train, y_train_encoded) # Use the encoded labels for training
8 y_pred_logreg = clf.predict(X_test)
9
10 # Train the XGBoost model
11 from xgboost import XGBClassifier
12 xgb_clf = XGBClassifier(use_label_encoder=False, eval_metric='mlogloss', random
13 xgb_clf.fit(X_train, y_train_encoded)
14 y_pred_xgb = xgb_clf.predict(X_test)
15
16 # Create confusion matrix plots
17 fig, axes = plt.subplots(1, 2, figsize=(14, 6))
18
19 # Logistic Regression Confusion Matrix
20 cm_logreg = confusion_matrix(y_test_encoded, y_pred_logreg)
21 ConfusionMatrixDisplay(confusion_matrix=cm_logreg, display_labels=clf.classes_)
22 axes[0].set_title("Logistic Regression")
23
24 # XGBoost Confusion Matrix
25 cm_xgb = confusion_matrix(y_test_encoded, y_pred_xgb)
26 ConfusionMatrixDisplay(confusion_matrix=cm_xgb, display_labels=xgb_clf.classes_)
27 axes[1].set_title("XGBoost")
28
29 plt.tight_layout()
30 plt.show()
31

```

➡ /usr/local/lib/python3.11/dist-packages/xgboost/core.py:158: UserWarning: [00: 1 Parameters: { "use_label_encoder" } are not used.

