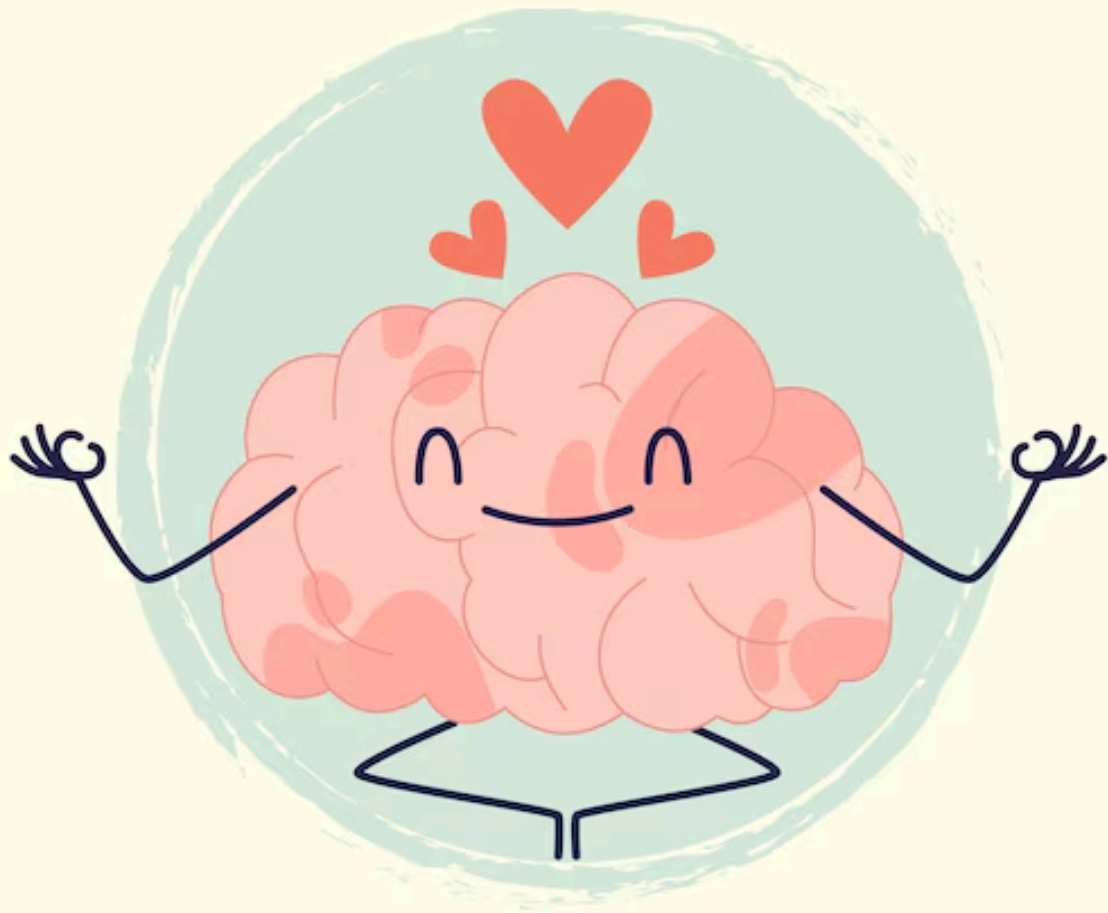


## ✓ Sentiment Analysis for Mental Health Monitoring 🤖

[dataset source](#) 📎



**MENTAL  
HEALTH  
awareness**

## About The Dataset

This comprehensive dataset is a meticulously curated collection of mental health statuses tagged from various statements. The dataset amalgamates raw data from multiple sources, cleaned and compiled to create a robust resource for developing chatbots and performing sentiment analysis.

---

## Data Source

The dataset integrates information from the following Kaggle datasets:

- 3k Conversations Dataset for Chatbot
  - Depression Reddit Cleaned
  - Human Stress Prediction
  - Predicting Anxiety in Mental Health Data
  - Mental Health Dataset Bipolar
  - Reddit Mental Health Data
  - Students Anxiety and Depression Dataset
  - Suicidal Mental Health Dataset
  - Suicidal Tweet Detection Dataset
- 

## Dataset Overview :

- **Description:**

This dataset is a comprehensive collection of 50,000 text statements related to mental health, each tagged with one of seven mental health statuses. The primary purpose of this dataset is to assist in building machine learning models for classifying mental health conditions based on textual data, such as social media posts or other user-generated content.

- **Columns:**

1. **unique\_id**: A unique identifier for each entry.
2. **statement**: A piece of text, typically a statement or comment, associated with a particular mental health status.
3. **status**: The mental health status assigned to the statement. The possible categories are:
  - Normal
  - Depression
  - Suicidal
  - Anxiety
  - Stress
  - Bi-Polar

## ■ Personality Disorder

---

### Usage :

This dataset is ideal for training machine learning models aimed at understanding and predicting mental health conditions based on textual data. It can be used in various applications such as:

- Chatbot development for mental health support.
  - Sentiment analysis to gauge mental health trends.
  -
- 

### How we use NLP Concepts

In our text processing methodology, we begin by removing punctuation, URLs, and hyperlinks from the statements. Additionally, we eliminate stop words, such as "is," "are," and "the," to enhance the focus on more informative terms.

In the realm of morphological analysis, we extract the stems of the words. Subsequently, we employ a TF-IDF (Term Frequency-Inverse Document Frequency) vectorizer to transform the processed text into a vector representation. This vector is then utilized to fit the mode

---

### Conlusion

in the conclusion we train and tune a four model wich is :

- SVM
- Logistic Regression
- Neural Ntwork
- KNN the first three give a Convergent results with accuracy near to 76 % , and the last one give a 65 % acc .

## ✓ Import Libraries :

lets start by importing the nessacary libraries

`pip install stanza`



Collecting stanza

Downloading stanza-1.9.2-py3-none-any.whl.metadata (13 kB)

Collecting emoji (from stanza)

Downloading emoji-2.14.0-py3-none-any.whl.metadata (5.7 kB)

Requirement already satisfied: numpy in /usr/local/lib/python3.10/dist-packages (from stanza)

Requirement already satisfied: protobuf>=3.15.0 in /usr/local/lib/python3.10/dist-packages (from stanza)

Requirement already satisfied: requests in /usr/local/lib/python3.10/dist-packages (from stanza)

Requirement already satisfied: networkx in /usr/local/lib/python3.10/dist-packages (from stanza)

Requirement already satisfied: torch>=1.3.0 in /usr/local/lib/python3.10/dist-packages (from stanza)

Requirement already satisfied: tqdm in /usr/local/lib/python3.10/dist-packages (from stanza)

Requirement already satisfied: tomli in /usr/local/lib/python3.10/dist-packages (from stanza)

Requirement already satisfied: filelock in /usr/local/lib/python3.10/dist-packages (from stanza)

Requirement already satisfied: typing-extensions>=4.8.0 in /usr/local/lib/python3.10/dist-packages (from stanza)

Requirement already satisfied: jinja2 in /usr/local/lib/python3.10/dist-packages (from stanza)

Requirement already satisfied: fsspec in /usr/local/lib/python3.10/dist-packages (from stanza)

Requirement already satisfied: sympy==1.13.1 in /usr/local/lib/python3.10/dist-packages (from stanza)

Requirement already satisfied: mpmath<1.4,>=1.1.0 in /usr/local/lib/python3.10/dist-packages (from stanza)

Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.10/dist-packages (from stanza)

Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.10/dist-packages (from stanza)

Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.10/dist-packages (from stanza)

Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.10/dist-packages (from stanza)

Requirement already satisfied: MarkupSafe>=2.0 in /usr/local/lib/python3.10/dist-packages (from stanza)

Downloading stanza-1.9.2-py3-none-any.whl (1.1 MB)

1.1/1.1 MB 42.4 MB/s eta 0:00:00

Downloading emoji-2.14.0-py3-none-any.whl (586 kB)

586.9/586.9 kB 24.1 MB/s eta 0:00:00

Installing collected packages: emoji, stanza

Successfully installed emoji-2.14.0 stanza-1.9.2

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline
from nltk.corpus import stopwords
from imblearn.over_sampling import SMOTE
```

```
import re
import random
#from imblearn.over_sampling import RandomOverSampler
from scipy.sparse import hstack # To combine sparse matrices
from wordcloud import WordCloud
from tensorflow.keras.models import Sequential
```

```
from tensorflow.keras.layers import Dense, Dropout
from tensorflow.keras.optimizers import Adam
from tensorflow.keras.layers import BatchNormalization
from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
from tensorflow.keras.regularizers import l2
```

```
from sklearn.svm import SVC
from sklearn.model_selection import GridSearchCV
from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
from sklearn.neighbors import KNeighborsClassifier
from sklearn.model_selection import GridSearchCV
from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import GridSearchCV
from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
```

```
import nltk
from nltk.tokenize import word_tokenize
from nltk.stem import PorterStemmer
```

```
import stanza
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.model_selection import train_test_split, GridSearchCV
from sklearn.preprocessing import LabelEncoder
from sklearn.linear_model import LogisticRegression
from sklearn.naive_bayes import BernoulliNB
from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
```

```
import warnings
warnings.filterwarnings("ignore")
```

```
import tensorflow as tf
```

```
# Check if GPU is available
print("Num GPUs Available: ", len(tf.config.list_physical_devices('GPU')))
```

 Num GPUs Available: 1

## › Datasert Loading and Preparing

[ ] ↳ 22 cells hidden

## ✓ NLP Pre-Processing

in this part of the ml pipeline, we perform several essential NLP preprocessing steps to prepare our data for machine learning. These steps include:

### 1. Text Processing:

- include :
  - pattern removing (links, emails and etc.)
  - lower casing
  - stop words removing

### 2. Tokenization & Stemming :

- Tokenization :
  - this step converts each text sample into an array of tokens (individual words or sub-words).
  - we use Stanza word tokenizer to ensure precise, language-specific tokenization.
- Stemming :
  - in this step, we transform each tokenized array into a stemmed version, reducing each word to its base or root form.
  - stanza NLP tools allow us to use various stemming techniques for extracting meaningful root words.

### 3. Part-of-Speech (POS) Tagging and Filtering:

- using Stanza, we apply POS tagging to identify and filter specific parts of speech, such as:
  - **Verbs**: Identify and include only verbs to capture action-related semantics.
  - **Nouns**: Extract nouns to focus on entity- or object-based information.
  - **Adjectives**: Include adjectives to analyze descriptive language and sentiment.
- This POS filtering allows us to tailor the input data by emphasizing different linguistic elements and can be customized based on the classification or NLP goals.

## > 1. Text Processing:

[ ] ↳ 10 cells hidden

## > 2. Tokenization & Stemming

[ ] ↳ 26 cells hidden

### ✓ 3. Part-Of-Speech (POS) Tagging

#### ✓ 3.1 Extract POS Taggs

##### > Nlp

[ ] ↳ 2 cells hidden

#### ✓ define function to extract pos

```
import nltk

def extract_pos_tokens(text, pos_tag_prefix):
    # Tokenize the input text
    tokens = nltk.word_tokenize(text)
    # Get the POS tags
    pos_tags = nltk.pos_tag(tokens)

    # Extract tokens based on POS tag prefix
    tokens_filtered = [word for word, pos in pos_tags if pos.startswith(pos_tag_prefix)]

    return ' '.join(tokens_filtered)
```

##### > Extract Verbs

[ ] ↳ 1 cell hidden

#### ✓ Extract Nouns

```
df['tokens_nouns'] = df['statement'].apply(lambda x: extract_pos_tokens(x, 'N'))
```

#### ✓ Extrac Adjectives

```
df['tokens_adj'] = df['statement'].apply(lambda x: extract_pos_tokens(x, 'J'))
```

#### ✓ Sample

```
df.sample(2)
```



	statement	status	statement_length	num_of_sentences	tokens	tokens_ste
24376	today left home mom took opportunity confront ...	Depression	502	1	[today, left, home, mom, took, opportunity, co...	today left h mom opp confront
	title says				title says	

## > About VVectorizing Approach

↳ 2 cells hidden

## > Data Spliting

[ ] ↳ 6 cells hidden

## ✓ Vectorizing

- Convert text to features using TF-IDF vectoriser

## > Vectorizing the stemmes

[ ] ↳ 1 cell hidden

## > Verbs-Vector

[ ] ↳ 1 cell hidden

## > Nouns-Vector

[ ] ↳ 2 cells hidden

## > Adjectives-Vector



[ ] ↳ 1 cell hidden

## ✓ Weighted-Vector

```
# Define weights for each POS feature type
verb_weight = 1.8
noun_weight = 3.0
adj_weight = 0.5

# Apply weights and combine features
X_weighted_train = hstack([
    X_verbs_train_tfidf * verb_weight,
    X_nouns_train_tfidf * noun_weight,
    X_adj_train_tfidf * adj_weight
])
# Apply weights and combine features
X_weighted_test = hstack([
    X_verbs_test_tfidf * verb_weight,
    X_nouns_test_tfidf * noun_weight,
    X_adj_test_tfidf * adj_weight
])
```

## › Resampling very importing remmember

[ ] ↳ 2 cells hidden

## ✓ Model Trainging

### › Prepare The Models

[ ] ↳ 4 cells hidden

### › Training Based on Verb-Vector

[ ] ↳ 16 cells hidden

### › Training Based on Nouns-Vector

[ ] ↴ 15 cells hidden

## ➤ Training Based on Adjective-Vector

[ ] ↴ 16 cells hidden

## ✓ Training Based on weighted-Vector

### ✓ ⌘ Logistic Regression Model

### ✓ Grid CV

```
# Fit the grid search on the training data
grid_search.fit(X_weighted_train, y_train)
```

➡ Fitting 5 folds for each of 8 candidates, totalling 40 fits

```
GridSearchCV ⓘ ?
  ▸ best_estimator_: LogisticRegression
    ▸ LogisticRegression ?
```

### ✓ Best params

```
# Get the best model and hyperparameters
best_clf = grid_search.best_estimator_
print("Best hyperparameters found: ", grid_search.best_params_)
```

➡ Best hyperparameters found: {'C': 1, 'penalty': 'l1', 'solver': 'liblinear'}

### ✓ Confusion Matrix

```
# Predict on the test set using the best model
y_pred = best_clf.predict(X_weighted_test)
```

```
# Calculate accuracy
accuracy_reg = accuracy_score(y_test, y_pred)
print("\nAccuracy: ", accuracy_reg)

# Compute the confusion matrix and classification report
conf_matrix_reg = confusion_matrix(y_test, y_pred)
labels = lbl_enc.classes_
print(classification_report(y_test, y_pred, target_names=labels))
```

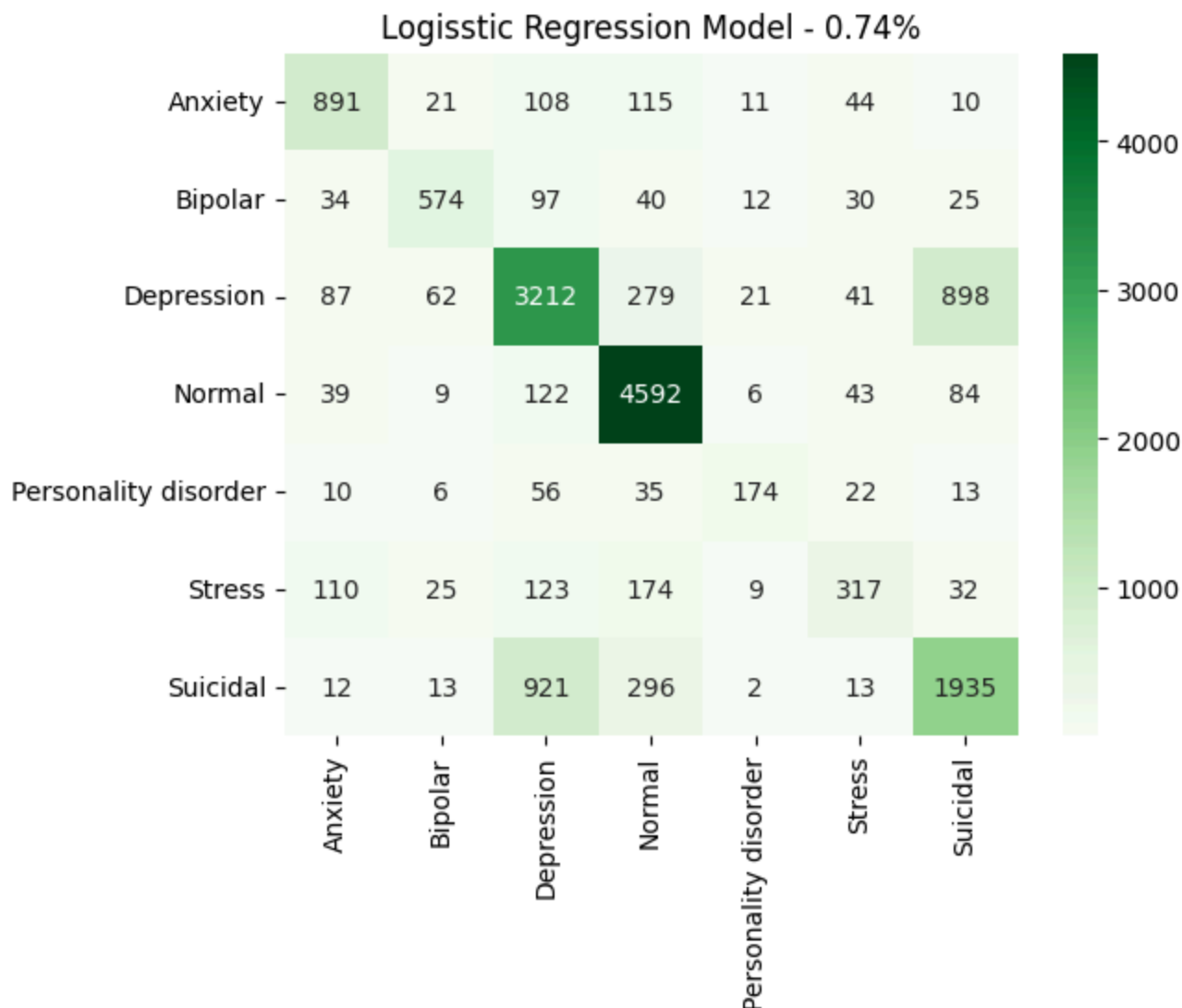


```
Accuracy: 0.7399557102182853
```

	precision	recall	f1-score	support
Anxiety	0.75	0.74	0.75	1200
Bipolar	0.81	0.71	0.75	812
Depression	0.69	0.70	0.70	4600
Normal	0.83	0.94	0.88	4895
Personality disorder	0.74	0.55	0.63	316
Stress	0.62	0.40	0.49	790
Suicidal	0.65	0.61	0.63	3192
accuracy			0.74	15805
macro avg	0.73	0.66	0.69	15805
weighted avg	0.73	0.74	0.73	15805

## ✓ Heat Map

```
ax = sns.heatmap(conf_matrix_reg, annot = True, fmt='d', cmap='Greens', xticklabels=labels,
ax.set_title(f'Logisstic Regression Model - {accuracy_reg:.2}%')
plt.show()
```



## ✓ Nerual Network

```
# Define the layers in an array
layers = [
    Dense(units=128, activation='relu', input_shape=(X_weighted_train.shape[1],), kernel_regularizer=BatchNormalization(),
    Dropout(rate=0.2), # Dropout Layer 1
    Dense(units=64, activation='relu', kernel_regularizer=l2(0.01)), # Hidden Layer 1
    Dropout(rate=0.1), # Dropout Layer 1
    Dense(units=16, activation='relu', kernel_regularizer=l2(0.01)), # Hidden Layer 2

    Dropout(rate=0.2), # Dropout Layer 2
    Dense(units=len(lbl_enc.classes_), activation='softmax') # Output Layer
]

# Initialize the Sequential model
model = Sequential(layers)
```

```
# Compile the model: Using Adam optimizer, sparse categorical crossentropy loss, and accuracy
model.compile(optimizer=Adam(learning_rate=0.001),
              loss='sparse_categorical_crossentropy',
              metrics=['accuracy'])
```

```
# Train the model on training data with validation on the test data
history = model.fit(X_weighted_train.toarray(),
                   y_train,
                   epochs=20,
                   batch_size=64,
                   validation_split=0.3
                   )
```

```
Epoch 1/20
404/404 ————— 12s 21ms/step - accuracy: 0.5236 - loss: 3.0971 - val_accu
Epoch 2/20
404/404 ————— 1s 3ms/step - accuracy: 0.7024 - loss: 1.3596 - val_accu
Epoch 3/20
404/404 ————— 1s 3ms/step - accuracy: 0.7218 - loss: 1.1369 - val_accu
Epoch 4/20
404/404 ————— 1s 3ms/step - accuracy: 0.7198 - loss: 1.1135 - val_accu
Epoch 5/20
404/404 ————— 1s 3ms/step - accuracy: 0.7320 - loss: 1.0978 - val_accu
Epoch 6/20
404/404 ————— 1s 3ms/step - accuracy: 0.7340 - loss: 1.0924 - val_accu
Epoch 7/20
404/404 ————— 1s 3ms/step - accuracy: 0.7462 - loss: 1.0642 - val_accu
Epoch 8/20
404/404 ————— 1s 3ms/step - accuracy: 0.7403 - loss: 1.0653 - val_accu
Epoch 9/20
404/404 ————— 1s 3ms/step - accuracy: 0.7440 - loss: 1.0474 - val_accu
Epoch 10/20
404/404 ————— 2s 4ms/step - accuracy: 0.7521 - loss: 1.0475 - val_accu
Epoch 11/20
404/404 ————— 3s 4ms/step - accuracy: 0.7546 - loss: 1.0314 - val_accu
Epoch 12/20
404/404 ————— 2s 3ms/step - accuracy: 0.7582 - loss: 1.0191 - val_accu
Epoch 13/20
404/404 ————— 1s 3ms/step - accuracy: 0.7661 - loss: 1.0047 - val_accu
Epoch 14/20
404/404 ————— 1s 3ms/step - accuracy: 0.7669 - loss: 0.9933 - val_accu
Epoch 15/20
404/404 ————— 1s 3ms/step - accuracy: 0.7709 - loss: 1.0023 - val_accu
Epoch 16/20
404/404 ————— 1s 3ms/step - accuracy: 0.7757 - loss: 0.9801 - val_accu
Epoch 17/20
404/404 ————— 1s 3ms/step - accuracy: 0.7697 - loss: 0.9858 - val_accu
Epoch 18/20
404/404 ————— 1s 3ms/step - accuracy: 0.7784 - loss: 0.9750 - val_accu
```

Epoch 19/20

**404/404** ————— 2s 4ms/step - accuracy: 0.7839 - loss: 0.9565 - val\_accu

Epoch 20/20

**404/404** ————— 2s 4ms/step - accuracy: 0.7887 - loss: 0.9545 - val\_accu

```

# Make predictions on the test set
y_pred_prob = model.predict(X_weighted_test.toarray())
y_pred = y_pred_prob.argmax(axis=1) # Convert probabilities to class predictions

# Calculate the accuracy
accuracy_nn = accuracy_score(y_test, y_pred)
print("\n")
print("Accuracy:", accuracy_nn)

# Compute the confusion matrix
labels = lbl_enc.classes_
conf_matrix_nn = confusion_matrix(y_test, y_pred)

# Print classification report
print(classification_report(y_test, y_pred, target_names=labels))

```

 **494/494** ————— 1s 2ms/step

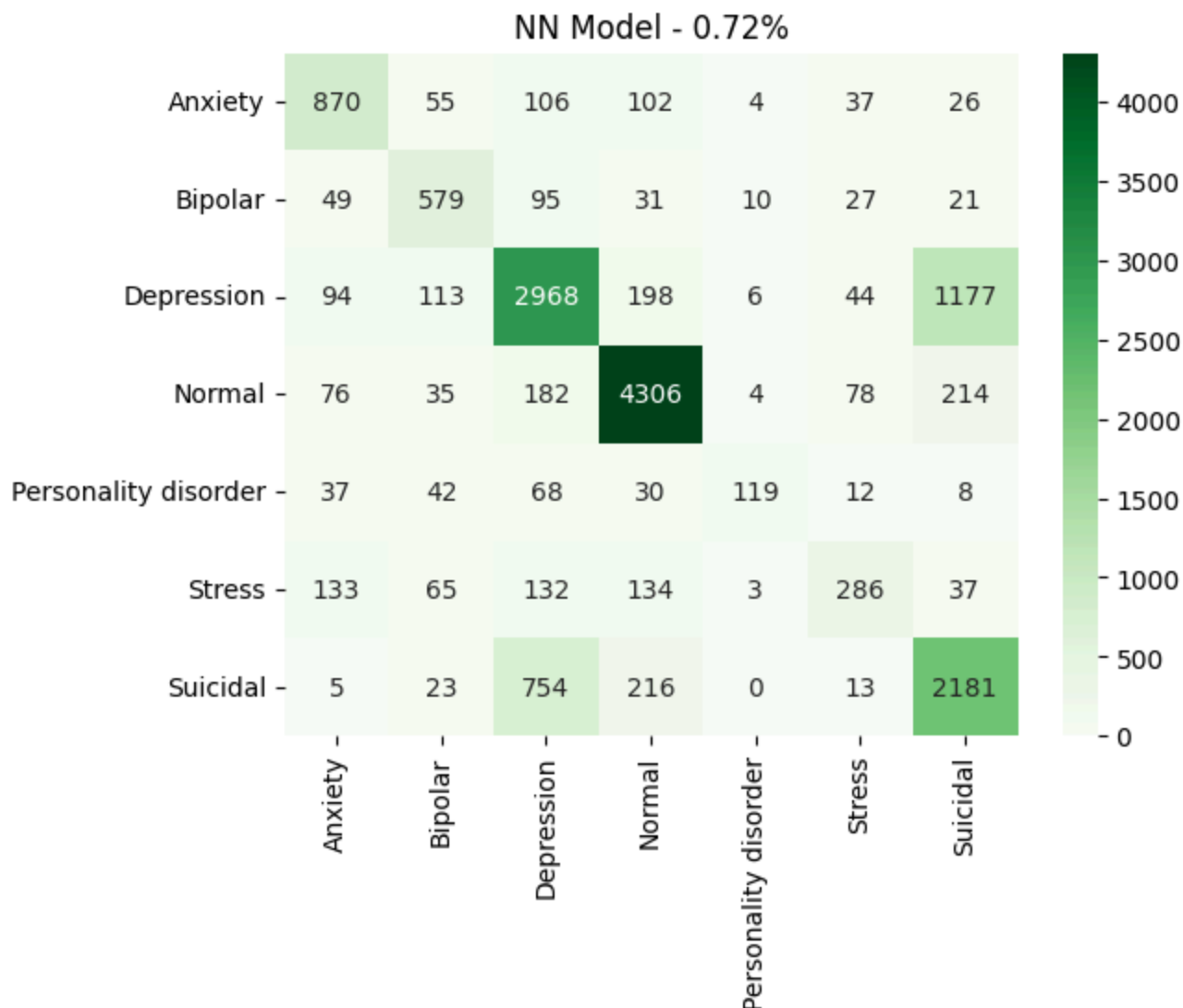
Accuracy: 0.7155330591584942

	precision	recall	f1-score	support
Anxiety	0.69	0.72	0.71	1200
Bipolar	0.63	0.71	0.67	812
Depression	0.69	0.65	0.67	4600
Normal	0.86	0.88	0.87	4895
Personality disorder	0.82	0.38	0.52	316
Stress	0.58	0.36	0.44	790
Suicidal	0.60	0.68	0.64	3192
accuracy			0.72	15805
macro avg	0.69	0.63	0.64	15805
weighted avg	0.72	0.72	0.71	15805

```

ax = sns.heatmap(conf_matrix_nn, annot = True, fmt='d', cmap='Greens', xticklabels=labels, y
ax.set_title(f'NN Model - {accuracy_nn:.2}%')
plt.show()

```



```
plt.figure(figsize=(12, 5))
```

```
# Plot training & validation loss values
```

```
plt.subplot(1, 2, 1) # 1 row, 2 columns, 1st subplot
```

```
plt.plot(history.history['loss'], label='Training Loss')
```

```
plt.plot(history.history['val_loss'], label='Validation Loss')
```

```
plt.title('Model Loss')
```

```
plt.xlabel('Epoch')
```

```
plt.ylabel('Loss')
```

```
plt.legend(loc='upper right')
```

```
# Plot training & validation accuracy values
```

```
plt.subplot(1, 2, 2) # 1 row, 2 columns, 2nd subplot
```

```
plt.plot(history.history['accuracy'], label='Training Accuracy')
```

```
plt.plot(history.history['val_accuracy'], label='Validation Accuracy')
```

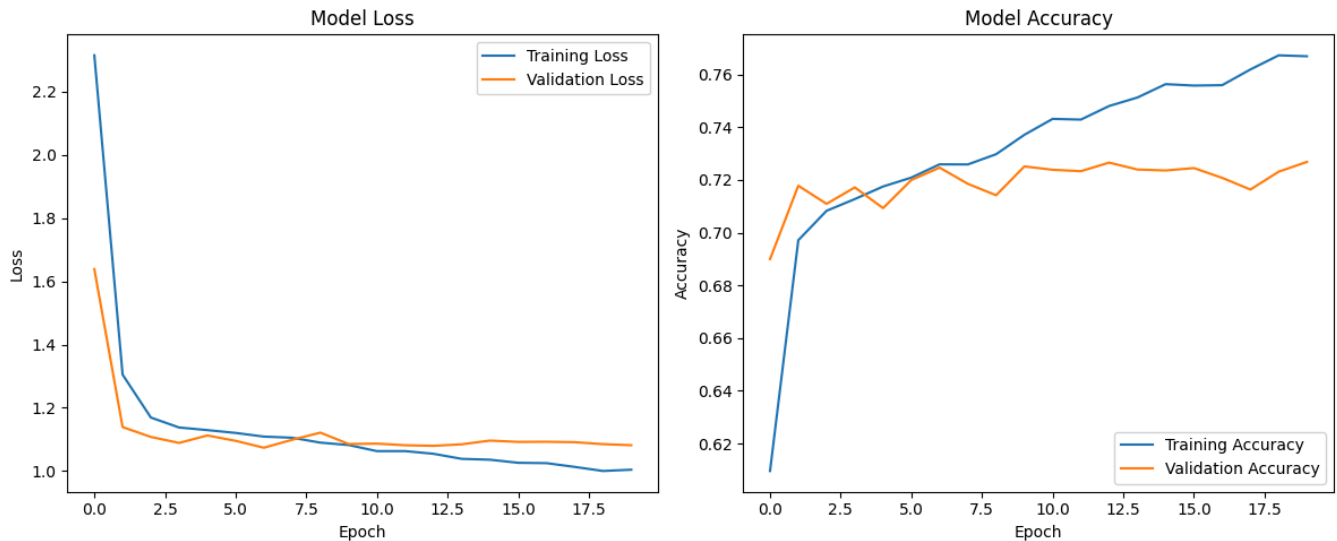
```
plt.title('Model Accuracy')
```

```
plt.xlabel('Epoch')
```

```
plt.ylabel('Accuracy')
```

```
plt.legend(loc='lower right')
```

```
plt.tight_layout() # Adjusts subplots to fit in the figure area.  
plt.show() # Display the plots
```



## › Training Based on stemm-Vector

[ ] ↳ 16 cells hidden

End



