

This notebook runs the gradient boosting machine with the best parameters (from gradient_boosting_machine.ipynb) using Faizan's preprocessing tools excluding the augmented text function using back translation

Accuracy: 0.752945612216043

```
In [1]: import pandas as pd
import numpy as np
import re
import string
import nltk
from nltk.corpus import stopwords
from nltk.tokenize import word_tokenize
from sklearn.model_selection import train_test_split, GridSearchCV
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.ensemble import GradientBoostingClassifier
from sklearn.metrics import classification_report, accuracy_score
```

```
In [2]: # Load the data
path = './kaggle_sentiment_data.csv'
df = pd.read_csv(path)
```

```
In [3]: # Handle NaN values in the statement column
df['statement'] = df['statement'].fillna('')
```

```
In [4]: # Data Preprocessing
def preprocess_text(text):
    text = text.lower() # Lowercase text
    text = re.sub(r'\[.*?\]', '', text) # Remove text in square brackets
    text = re.sub(r'https?://\S+|www\.\S+', '', text) # Remove links
    text = re.sub(r'<.*?>+', '', text) # Remove HTML tags
    text = re.sub(r'[%s]' % re.escape(string.punctuation), '', text) # Remove punctuation
    text = re.sub(r'\n', '', text) # Remove newlines
    text = re.sub(r'\w*\d\w*', '', text) # Remove words containing numbers
    return text
```

```
In [5]: # Tokenization and Stopwords Removal
stop_words = set(stopwords.words('english'))

def remove_stopwords(text):
    tokens = word_tokenize(text)
    tokens = [word for word in tokens if word not in stop_words]
    return ' '.join(tokens)
```

```
In [6]: # Preprocess the text data
df['cleaned_statement'] = df['statement'].apply(preprocess_text).apply(remove_stopwords)

# Ensure no NaN values
df['cleaned_statement'] = df['cleaned_statement'].fillna('')

# Splitting the data
X = df['cleaned_statement']
```

```

y = df['status']

# Vectorization
vectorizer = TfidfVectorizer(max_features=10000)
X_vec = vectorizer.fit_transform(X)

X_train, X_test, y_train, y_test = train_test_split(X_vec, y, test_size=0.2,

```

```

In [7]: # Gradient Boosting Model
gbm = GradientBoostingClassifier(learning_rate=0.1, max_depth=5, n_estimators=300)

```

```

In [8]: gbm.fit(X_train, y_train)

```

```

Out[8]: ▼ GradientBoostingClassifier
GradientBoostingClassifier(max_depth=5, n_estimators=300)

```

```

In [9]: y_pred = gbm.predict(X_test)

# Evaluate performance
print("Accuracy:", accuracy_score(y_test, y_pred))
print("Classification Report:\n", classification_report(y_test, y_pred))

```

Accuracy: 0.752945612216043

Classification Report:

	precision	recall	f1-score	support
Anxiety	0.81	0.72	0.76	778
Bipolar	0.85	0.70	0.77	575
Depression	0.69	0.73	0.71	3081
Normal	0.81	0.95	0.88	3270
Personality disorder	0.80	0.46	0.58	240
Stress	0.67	0.50	0.57	534
Suicidal	0.70	0.61	0.65	2131
accuracy			0.75	10609
macro avg	0.76	0.67	0.70	10609
weighted avg	0.75	0.75	0.75	10609