## ‣ Datasets

```
[ ] ↳ 9 cells hidden
```

## ▾ Classification Modeling on Sentiment Prediction

```
1 # Create a copy of the bitcoin price DataFrame
2 crypto_usd.head(2)
```

|   | time | close | high | low | open | volumefrom | volumeto |
|---|------|-------|------|-----|------|------------|----------|
| 0 | 2023-02-19 13:00:00 | 24682.03 | 24715.82 | 24682.03 | 24707.39 | 903.97 | 22335943.28 |
|   | 2023- | | | | | | |

```
1 print(crypto_usd.columns)

    Index(['time', 'close', 'high', 'low', 'open', 'volumefrom', 'volumeto',
           'Date', 'Time', 'volume', 'marketcap', 'price_delta'],
          dtype='object')
```

```
1 # Create a copy of the  bitcoin tweets DataFrame
2 df_tweets = tweets.copy()
3 df_tweets.head(2)
```

|   | user_name | user_location | user_description | user_created | user_follower |
|---|-----------|---------------|------------------|--------------|---------------|
| 0 | Irk | Vancouver, WA | Irk started investing in the stock market in 1... | 2018-08-11 03:17:00 | 116 |
| 1 | Xiang Zhang | NaN | Professional Software Engineer ð□□»ð□□□Crypto ... | 2011-01-11 01:37:00 | 42 |

```
1 # Merge the tweet data with the Bitcoin price data
2 tweets_df = pd.merge(df_tweets, crypto_usd, left_on='date', right_on='time', how='inner')
```

```
1 print(tweets_df.columns)
2

    Index(['user_name', 'user_location', 'user_description', 'user_created',
           'user_followers', 'user_friends', 'user_favourites', 'user_verified',
           'date', 'text', 'hashtags', 'source', 'is_retweet', 'compound', 'score',
           'sentiment_level', 'polarity', 'subjectivity', 'time', 'close', 'high',
           'low', 'open', 'volumefrom', 'volumeto', 'Date', 'Time', 'volume',
           'marketcap', 'price_delta'],
          dtype='object')
```

## ▾ Feature Extraction

```
1 import pandas as pd
2 from sklearn.feature_extraction.text import CountVectorizer
3 from sklearn.model_selection import train_test_split
4 from sklearn.naive_bayes import MultinomialNB
```

```
5 from sklearn.metrics import accuracy_score
6 from scipy.sparse import hstack
7
8 # Feature Extraction: Unigrams
9 unigram_vectorizer = CountVectorizer(ngram_range=(1, 1))
10 unigram_features = unigram_vectorizer.fit_transform(tweets_df['text'])
11
12 # Feature Extraction: Bigrams
13 bigram_vectorizer = CountVectorizer(ngram_range=(2, 2))
14 bigram_features = bigram_vectorizer.fit_transform(tweets_df['text'])
15
16 # Combining Features
17 combined_features = hstack([unigram_features, bigram_features])
18
19 # Perform sentiment analysis
20 X = combined_features
21 y = tweets_df['sentiment_level']
22
23 # Split the data into training and testing sets
24 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

```
1 #from sklearn.feature_extraction.text import CountVectorizer: This line imports the CountVectorizer class from the Scikit-learn library.
```

```
1 import numpy as np
2
3 # Print the first 10 rows of the term frequency matrix
4 print(combined_features[:10].toarray())
5
```

```
[[0 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]
 ...
 [0 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]]
```

## ▾ Naive_bayes

```
1 from sklearn.metrics import classification_report
```

```
1 from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
2
3 # Train a classification model (e.g., Naive Bayes)
4 classifier = MultinomialNB()
5 classifier.fit(X_train, y_train)
6
7 # Predict sentiment labels for test data
8 y_pred = classifier.predict(X_test)
9
10 # Evaluate the model using additional metrics
11 accuracy = accuracy_score(y_test, y_pred)
12 precision = precision_score(y_test, y_pred, average='weighted')
13 recall = recall_score(y_test, y_pred, average='weighted')
14 f1 = f1_score(y_test, y_pred, average='weighted')
15
16 print("Accuracy:", accuracy)
17 print("Precision:", precision)
18 print("Recall:", recall)
19 print("F1-Score:", f1)
20
21 # Use the trained model for future predictions
22 new_tweet = ["New tweet about Bitcoin"]
23 new_tweet_features = hstack([unigram_vectorizer.transform(new_tweet), bigram_vectorizer.transform(new_tweet)])
24 predicted_sentiment = classifier.predict(new_tweet_features)
25 #Classification Report
26 print("Predicted sentiment:", predicted_sentiment)
27 print(classification_report(y_test, y_pred))
```

```
Accuracy: 0.7879746835443038
Precision: 0.8031951087547753
Recall: 0.7879746835443038
F1-Score: 0.791715079259514
```

```
Predicted sentiment: ['Neutral']
                 precision    recall  f1-score   support

Extreme Negative      0.93      0.71      0.80        55
Extreme Positive      0.50      0.69      0.58       127
        Negative      0.83      0.61      0.71       157
         Neutral      0.88      0.86      0.87       908
        Positive      0.67      0.74      0.70       333

        accuracy                          0.79      1580
       macro avg      0.76      0.72      0.73      1580
    weighted avg      0.80      0.79      0.79      1580
```

```
1 #The Naive Bayes classifier achieved an accuracy of 0.7879746835443038 and a precision of 0.8031951087547753 for sentiment analysis on th
```

## ▾ Support Vector Machines (SVM)

```
 1 import pandas as pd
 2 from sklearn.feature_extraction.text import CountVectorizer
 3 from sklearn.model_selection import train_test_split
 4 from sklearn.svm import LinearSVC
 5 from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
 6
 7 try:
 8     # Split the data into training and testing sets
 9     X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
10
11     # Train a linear SVM classifier
12     classifier = LinearSVC()
13     classifier.fit(X_train, y_train)
14
15     # Evaluate the model using additional metrics
16     y_pred = classifier.predict(X_test)
17     accuracy = accuracy_score(y_test, y_pred)
18     precision = precision_score(y_test, y_pred, average='weighted')
19     recall = recall_score(y_test, y_pred, average='weighted')
20     f1 = f1_score(y_test, y_pred, average='weighted')
21
22     print("Accuracy:", accuracy)
23     print("Precision:", precision)
24     print("Recall:", recall)
25     print("F1-Score:", f1)
26
27     # Use the trained model for future predictions
28     new_tweet = ["New tweet about Bitcoin"]
29     new_tweet_features = hstack([unigram_vectorizer.transform(new_tweet), bigram_vectorizer.transform(new_tweet)])
30     predicted_sentiment = classifier.predict(new_tweet_features)
31     print("Predicted sentiment:", predicted_sentiment)
32
33 except Exception as e:
34     print("An error occurred:", str(e))
35 #Classification Report
36 print(classification_report(y_test, y_pred))
```

```
Accuracy: 0.8645569620253165
Precision: 0.8642113824767497
Recall: 0.8645569620253165
F1-Score: 0.8596632616414193
Predicted sentiment: ['Neutral']
                 precision    recall  f1-score   support

Extreme Negative      0.95      0.76      0.85        55
Extreme Positive      0.86      0.65      0.74       127
        Negative      0.89      0.70      0.78       157
         Neutral      0.87      0.97      0.92       908
        Positive      0.82      0.74      0.78       333

        accuracy                          0.86      1580
       macro avg      0.88      0.77      0.81      1580
    weighted avg      0.86      0.86      0.86      1580
```

```
1 #The SVM classifier achieved an accuracy of 0.8645569620253165 and a precision of 0.8642113824767497 for sentiment analysis on the tweet
```

## Random Forest

```python
1  import pandas as pd
2  from sklearn.feature_extraction.text import CountVectorizer
3  from sklearn.model_selection import train_test_split
4  from sklearn.ensemble import RandomForestClassifier
5  from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
6
7
8
9  # Feature Extraction: Unigrams
10 vectorizer = CountVectorizer(ngram_range=(1, 1))
11 X = vectorizer.fit_transform(tweets_df['text'])
12 y = tweets_df['sentiment_level']
13
14 try:
15     # Split the data into training and testing sets
16     X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
17
18     # Train a linear Random Forest classifier
19     classifier = RandomForestClassifier()
20     classifier.fit(X_train, y_train)
21
22     # Evaluate the model using additional metrics
23     y_pred = classifier.predict(X_test)
24     accuracy = accuracy_score(y_test, y_pred)
25     precision = precision_score(y_test, y_pred, average='weighted')
26     recall = recall_score(y_test, y_pred, average='weighted')
27     f1 = f1_score(y_test, y_pred, average='weighted')
28
29     print("Accuracy:", accuracy)
30     print("Precision:", precision)
31     print("Recall:", recall)
32     print("F1-Score:", f1)
33
34     # Use the trained model for future predictions
35     new_tweet = ["New tweet about Bitcoin"]
36     new_tweet_features = vectorizer.transform(new_tweet)
37     predicted_sentiment = classifier.predict(new_tweet_features)
38     print("Predicted sentiment:", predicted_sentiment)
39
40 except Exception as e:
41     print("An error occurred:", str(e))
42 #Classification Report
43 print(classification_report(y_test, y_pred))
```

```
Accuracy: 0.8575949367088608
Precision: 0.8647362580034188
Recall: 0.8575949367088608
F1-Score: 0.8488711984107943
Predicted sentiment: ['Neutral']
                  precision    recall  f1-score   support

Extreme Negative       1.00      0.75      0.85        55
Extreme Positive       0.94      0.47      0.63       127
        Negative       0.95      0.66      0.78       157
         Neutral       0.85      0.99      0.91       908
        Positive       0.81      0.77      0.79       333

        accuracy                           0.86      1580
       macro avg       0.91      0.73      0.79      1580
    weighted avg       0.86      0.86      0.85      1580
```

```python
1  #The Random Forest classifier achieved an accuracy of 0.870253164556962 and a precision of 0.8670628029958947 for sentiment analysis on t
```

## Logistic Regression

```python
1  import pandas as pd
2  from sklearn.feature_extraction.text import CountVectorizer
3  from sklearn.model_selection import train_test_split
4  from sklearn.linear_model import LogisticRegression
```

```
 5 from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
 6
 7 # Assuming you have tweets_df with the appropriate 'text' and 'sentiment_level' columns
 8
 9 # Feature Extraction: Unigrams
10 vectorizer = CountVectorizer(ngram_range=(1, 1))
11 X = vectorizer.fit_transform(tweets_df['text'])
12 y = tweets_df['sentiment_level']
13
14 try:
15     # Split the data into training and testing sets
16     X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
17
18     # Train a logistic regression classifier with increased max_iter
19     classifier = LogisticRegression(max_iter=1000)
20     classifier.fit(X_train, y_train)
21
22     # Evaluate the model using additional metrics
23     y_pred = classifier.predict(X_test)
24     accuracy = accuracy_score(y_test, y_pred)
25     precision = precision_score(y_test, y_pred, average='weighted')
26     recall = recall_score(y_test, y_pred, average='weighted')
27     f1 = f1_score(y_test, y_pred, average='weighted')
28
29     print("Accuracy:", accuracy)
30     print("Precision:", precision)
31     print("Recall:", recall)
32     print("F1-Score:", f1)
33
34     # Use the trained model for future predictions
35     new_tweet = ["New tweet about Bitcoin"]
36     new_tweet_features = vectorizer.transform(new_tweet)
37     predicted_sentiment = classifier.predict(new_tweet_features)
38     print("Predicted sentiment:", predicted_sentiment)
39
40 except Exception as e:
41     print("An error occurred:", str(e))
42 #Classification Report
43 print(classification_report(y_test, y_pred))
```

```
    Accuracy: 0.859493670886076
    Precision: 0.8596313804881421
    Recall: 0.859493670886076
    F1-Score: 0.8545443425051455
    Predicted sentiment: ['Neutral']
                    precision    recall  f1-score   support

    Extreme Negative     1.00      0.73      0.84        55
    Extreme Positive     0.86      0.64      0.73       127
            Negative     0.86      0.69      0.76       157
             Neutral     0.87      0.97      0.91       908
            Positive     0.82      0.75      0.78       333

            accuracy                         0.86      1580
           macro avg     0.88      0.75      0.81      1580
        weighted avg     0.86      0.86      0.85      1580
```

```
 1 #The Logistic Regression classifier achieved an accuracy of 0.859493670886076 and a precision of 0.8596313804881421 for sentiment analysi
```

## ▾ Gradient Boosting

```
 1 import pandas as pd
 2 from sklearn.feature_extraction.text import CountVectorizer
 3 from sklearn.model_selection import train_test_split
 4 from sklearn.ensemble import GradientBoostingClassifier
 5 from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
 6
 7 # Assuming you have tweets_df with the appropriate 'text' and 'sentiment_level' columns
 8
 9 # Feature Extraction: Unigrams
10 vectorizer = CountVectorizer(ngram_range=(1, 1))
11 X = vectorizer.fit_transform(tweets_df['text'])
12 y = tweets_df['sentiment_level']
13
```

```
14 try:
15     # Split the data into training and testing sets
16     X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
17
18     # Train a Gradient Boosting classifier
19     classifier = GradientBoostingClassifier()
20     classifier.fit(X_train, y_train)
21
22     # Evaluate the model using additional metrics
23     y_pred = classifier.predict(X_test)
24     accuracy = accuracy_score(y_test, y_pred)
25     precision = precision_score(y_test, y_pred, average='weighted')
26     recall = recall_score(y_test, y_pred, average='weighted')
27     f1 = f1_score(y_test, y_pred, average='weighted')
28
29     print("Accuracy:", accuracy)
30     print("Precision:", precision)
31     print("Recall:", recall)
32     print("F1-Score:", f1)
33
34     # Use the trained model for future predictions
35     new_tweet = ["New tweet about Bitcoin"]
36     new_tweet_features = vectorizer.transform(new_tweet)
37     predicted_sentiment = classifier.predict(new_tweet_features)
38     print("Predicted sentiment:", predicted_sentiment)
39
40 except Exception as e:
41     print("An error occurred:", str(e))
42 #Classification Report
43 print(classification_report(y_test, y_pred))
```

```
   Accuracy: 0.8468354430379746
   Precision: 0.8543684218834472
   Recall: 0.8468354430379746
   F1-Score: 0.8375792142254445
   Predicted sentiment: ['Neutral']
                     precision    recall  f1-score   support

   Extreme Negative       0.90      0.78      0.83        55
   Extreme Positive       0.94      0.57      0.71       127
           Negative       0.92      0.62      0.74       157
            Neutral       0.83      0.99      0.90       908
           Positive       0.86      0.68      0.76       333

           accuracy                           0.85      1580
          macro avg       0.89      0.73      0.79      1580
       weighted avg       0.85      0.85      0.84      1580
```

```
1 #The Gradient Boosting classifier achieved an accuracy of 0.8468354430379746 and a precision of 0.8543684218834472 for sentiment analysis
```

## Cross Validation of Models

```
1 from sklearn.naive_bayes import MultinomialNB
2 from sklearn.svm import SVC
3 from sklearn.ensemble import RandomForestClassifier, GradientBoostingClassifier
4 from sklearn.linear_model import LogisticRegression
5 from sklearn.model_selection import cross_val_score
6
7 # Define the models
8 models = [
9     ("Naive Bayes", MultinomialNB()),
10     ("Support Vector Machine", SVC()),
11     ("Random Forest", RandomForestClassifier()),
12     ("Logistic Regression", LogisticRegression()),
13     ("Gradient Boosting", GradientBoostingClassifier())
14 ]
15
16 # Perform cross-validation and evaluation for each model
17 for model_name, model in models:
18     # Perform cross-validation
19     scores = cross_val_score(model, X_train, y_train, cv=5)
20     mean_score = scores.mean()
21
```

```
22      # Fit the model on the entire training set
23      model.fit(X_train, y_train)
24
25      # Evaluate the model on the test set
26      accuracy = model.score(X_test, y_test)
27
28      # Print the results
29      print("Model:", model_name)
30      print("Cross-Validation Mean Score:", mean_score)
31      print("Accuracy:", accuracy)
32      print()
33
```

```
Model: Naive Bayes
Cross-Validation Mean Score: 0.7910727171592651
Accuracy: 0.7879746835443038
```

## ▾ Cross Validation

```
1 #cross-validation for the models using scikit-learn's cross_val_score function
```

```
1 import pandas as pd
2 from sklearn.feature_extraction.text import CountVectorizer
3 from sklearn.feature_selection import SelectKBest, chi2
4 from sklearn.model_selection import train_test_split, cross_val_score
5 from sklearn.naive_bayes import MultinomialNB
6 from sklearn.svm import SVC
7 from sklearn.ensemble import RandomForestClassifier
8 from sklearn.metrics import accuracy_score
9 from scipy.sparse import hstack
10
```

```
1 #Naive Bayes
```

```
1 from sklearn.naive_bayes import MultinomialNB
2 from sklearn.svm import LinearSVC
3 from sklearn.ensemble import RandomForestClassifier
4 from sklearn.model_selection import cross_val_score
5
6 # Train and evaluate Naive Bayes
7 naive_bayes = MultinomialNB()
8 naive_bayes_scores = cross_val_score(naive_bayes, X_train, y_train, cv=5)
9 print("Naive Bayes Cross-Validation Scores:", naive_bayes_scores.mean())
10 naive_bayes.fit(X_train, y_train)
11 naive_bayes_accuracy = naive_bayes.score(X_test, y_test)
12 print("Naive Bayes Accuracy:", naive_bayes_accuracy)
13 # Predict sentiment labels for test data
14 y_pred = naive_bayes.predict(X_test)
15 from sklearn.metrics import classification_report
16 print(classification_report(y_test, y_pred))
```

```
Naive Bayes Cross-Validation Scores: 0.7888584042414585
Naive Bayes Accuracy: 0.7943037974683544
                   precision    recall  f1-score   support

Extreme Negative       0.94      0.58      0.72        55
Extreme Positive       0.60      0.57      0.59       127
        Negative       0.85      0.58      0.69       157
         Neutral       0.84      0.92      0.88       908
        Positive       0.69      0.68      0.69       333

        accuracy                           0.79      1580
       macro avg       0.79      0.67      0.71      1580
    weighted avg       0.79      0.79      0.79      1580
```

```
1 #SVM
```

```
1 from sklearn.naive_bayes import MultinomialNB
2 from sklearn.svm import LinearSVC
3 from sklearn.ensemble import RandomForestClassifier
```

```
 4 from sklearn.model_selection import cross_val_score
 5
 6 # Train and evaluate SVM
 7 svm = LinearSVC()
 8 svm_scores = cross_val_score(svm, X_train, y_train, cv=5)
 9 print("SVM Cross-Validation Scores:", svm_scores.mean())
10 svm.fit(X_train, y_train)
11 svm_accuracy = svm.score(X_test, y_test)
12 print("SVM Accuracy:", svm_accuracy)
13 # Predict sentiment labels for test data
14 y_pred = svm.predict(X_test)
15 from sklearn.metrics import classification_report
16 print(classification_report(y_test, y_pred))
```

```
    SVM Cross-Validation Scores: 0.8630914439199415
    SVM Accuracy: 0.870253164556962
                    precision    recall  f1-score   support

    Extreme Negative     0.89      0.76      0.82        55
    Extreme Positive     0.79      0.70      0.74       127
            Negative     0.81      0.71      0.76       157
             Neutral     0.90      0.96      0.93       908
            Positive     0.83      0.77      0.80       333

            accuracy                         0.87      1580
           macro avg     0.84      0.78      0.81      1580
        weighted avg     0.87      0.87      0.87      1580
```

```
 1 #Random Forest
```

```
 1
 2 # Train Random Forest classifier
 3 random_forest = RandomForestClassifier(n_estimators=100, n_jobs=-1)
 4 random_forest.fit(X_train, y_train)
 5
 6 # Evaluate Random Forest
 7 random_forest_scores = cross_val_score(random_forest, X_train, y_train, cv=5)
 8 random_forest_mean_score = random_forest_scores.mean()
 9
10 random_forest_accuracy = random_forest.score(X_test, y_test)
11
12 # Print results
13 print("Random Forest Cross-Validation Mean Score:", random_forest_mean_score)
14 print("Random Forest Accuracy:", random_forest_accuracy)
15 # Predict sentiment labels for test data
16 y_pred = random_forest.predict(X_test)
17 from sklearn.metrics import classification_report
18 print(classification_report(y_test, y_pred))
19
```

```
    Random Forest Cross-Validation Mean Score: 0.8553332681880594
    Random Forest Accuracy: 0.8645569620253165
                    precision    recall  f1-score   support

    Extreme Negative     1.00      0.75      0.85        55
    Extreme Positive     0.96      0.52      0.67       127
            Negative     0.94      0.66      0.78       157
             Neutral     0.85      0.99      0.92       908
            Positive     0.83      0.78      0.80       333

            accuracy                         0.86      1580
           macro avg     0.92      0.74      0.80      1580
        weighted avg     0.87      0.86      0.86      1580
```

```
 1 #Logistic Regression
```

```
 1 from sklearn.naive_bayes import MultinomialNB
 2 from sklearn.svm import LinearSVC
 3 from sklearn.ensemble import RandomForestClassifier, GradientBoostingClassifier
 4 from sklearn.linear_model import LogisticRegression
 5 from sklearn.model_selection import cross_val_score
 6
 7 # Train and evaluate Logistic Regression
 8 logistic_regression = LogisticRegression(max_iter=1000)
```

```
 9 logistic_regression_scores = cross_val_score(logistic_regression, X_train, y_train, cv=5)
10 logistic_regression_mean_score = logistic_regression_scores.mean()
11 logistic_regression.fit(X_train, y_train)
12 logistic_regression_accuracy = logistic_regression.score(X_test, y_test)
13 print("Logistic Regression Cross-Validation Mean Score:", logistic_regression_mean_score)
14 print("Logistic Regression Accuracy:", logistic_regression_accuracy)
15 # Predict sentiment labels for test data
16 y_pred = logistic_regression.predict(X_test)
17 from sklearn.metrics import classification_report
18 print(classification_report(y_test, y_pred))
```

```
    Logistic Regression Cross-Validation Mean Score: 0.8429882387724625
    Logistic Regression Accuracy: 0.8537974683544304
                   precision    recall  f1-score   support

Extreme Negative       0.98      0.75      0.85        55
Extreme Positive       0.90      0.55      0.68       127
        Negative       0.92      0.69      0.79       157
         Neutral       0.84      0.98      0.91       908
        Positive       0.83      0.73      0.78       333

        accuracy                           0.85      1580
       macro avg       0.89      0.74      0.80      1580
    weighted avg       0.86      0.85      0.85      1580
```

```
1 #Gradient Boosting
```

```
 1 from sklearn.ensemble import GradientBoostingClassifier
 2 from sklearn.model_selection import cross_val_score
 3
 4 # Train and evaluate Gradient Boosting Classifier
 5 gradient_boosting = GradientBoostingClassifier()
 6 gradient_boosting_scores = cross_val_score(gradient_boosting, X_train, y_train, cv=3)  # Adjust cv parameter as needed
 7 gradient_boosting_mean_score = gradient_boosting_scores.mean()
 8
 9 gradient_boosting.fit(X_train, y_train)
10 gradient_boosting_accuracy = gradient_boosting.score(X_test, y_test)
11
12 print("Gradient Boosting Cross-Validation Mean Score:", gradient_boosting_mean_score)
13 print("Gradient Boosting Accuracy:", gradient_boosting_accuracy)
14 # Predict sentiment labels for test data
15 y_pred = gradient_boosting.predict(X_test)
16 from sklearn.metrics import classification_report
17 print(classification_report(y_test, y_pred))
```

```
    Gradient Boosting Cross-Validation Mean Score: 0.8421968977524533
    Gradient Boosting Accuracy: 0.8436708860759494
                   precision    recall  f1-score   support

Extreme Negative       0.89      0.76      0.82        55
Extreme Positive       0.93      0.58      0.71       127
        Negative       0.93      0.61      0.74       157
         Neutral       0.82      0.99      0.90       908
        Positive       0.86      0.66      0.75       333

        accuracy                           0.84      1580
       macro avg       0.89      0.72      0.78      1580
    weighted avg       0.85      0.84      0.83      1580
```

## ▾ Hyperparameter Tuning

```
 1 import pandas as pd
 2 from sklearn.feature_extraction.text import CountVectorizer
 3 from sklearn.feature_selection import SelectKBest, chi2
 4 from sklearn.model_selection import train_test_split, cross_val_score, GridSearchCV
 5 from sklearn.naive_bayes import MultinomialNB
 6 from sklearn.svm import SVC
 7 from sklearn.ensemble import RandomForestClassifier, GradientBoostingClassifier
 8 from sklearn.linear_model import LogisticRegression
 9 from sklearn.metrics import accuracy_score
10 from scipy.sparse import hstack
11
12 # Feature Extraction: Unigrams
```

```
13 unigram_vectorizer = CountVectorizer(ngram_range=(1, 1))
14 unigram_features = unigram_vectorizer.fit_transform(tweets_df['text'])
15
16 # Feature Extraction: Bigrams
17 bigram_vectorizer = CountVectorizer(ngram_range=(2, 2))
18 bigram_features = bigram_vectorizer.fit_transform(tweets_df['text'])
19
20 # Combining Features
21 combined_features = hstack([unigram_features, bigram_features])
22
23 # Perform sentiment analysis
24 X = combined_features
25 y = tweets_df['sentiment_level']
26
27 # Apply feature selection
28 k = 1000  # Number of top features to select
29 feature_selector = SelectKBest(chi2, k=k)
30 X_selected = feature_selector.fit_transform(X, y)
31
32 # Split the data into training and testing sets
33 X_train, X_test, y_train, y_test = train_test_split(X_selected, y, test_size=0.2, random_state=42)
34
35 # Define the models and their respective hyperparameter grids
36 models = [
37     ("Naive Bayes", MultinomialNB(), {'alpha': [0.1, 1.0, 10.0]}),
38     ("Support Vector Machine", SVC(), {'C': [0.1, 1.0, 10.0]}),
39     ("Random Forest", RandomForestClassifier(), {'n_estimators': [100, 200, 300]}),
40     ("Logistic Regression", LogisticRegression(), {'C': [0.1, 1.0, 10.0]}),
41     ("Gradient Boosting", GradientBoostingClassifier(), {'n_estimators': [100, 200, 300]})
42 ]
43
44 # Perform cross-validation and evaluation for each model
45 for model_name, model, param_grid in models:
46     # Perform hyperparameter tuning using GridSearchCV
47     grid_search = GridSearchCV(model, param_grid, cv=5)
48     grid_search.fit(X_train, y_train)
49
50     # Get the best model and its parameters
51     best_model = grid_search.best_estimator_
52     best_params = grid_search.best_params_
53
54     # Perform cross-validation with the best model
55     cross_val_scores = cross_val_score(best_model, X_train, y_train, cv=5)
56
57     # Fit the best model on the entire training set
58     best_model.fit(X_train, y_train)
59
60     # Make predictions on the test set
61     y_pred = best_model.predict(X_test)
62
63     # Calculate accuracy
64     accuracy = accuracy_score(y_test, y_pred)
65
66     # Print the results
67     print("Model:", model_name)
68     print("Best Parameters:", best_params)
69     print("Cross-Validation Accuracy:", cross_val_scores.mean())
70     print("Accuracy:", accuracy)
71     print()
72
```

```
Model: Naive Bayes
Best Parameters: {'alpha': 1.0}
Cross-Validation Accuracy: 0.7736638954869359
Accuracy: 0.7639240506329114

Model: Support Vector Machine
Best Parameters: {'C': 10.0}
Cross-Validation Accuracy: 0.8809735710634715
Accuracy: 0.8848101265822785

Model: Random Forest
Best Parameters: {'n_estimators': 300}
Cross-Validation Accuracy: 0.8649886747446806
Accuracy: 0.8613924050632912

/usr/local/lib/python3.10/dist-packages/sklearn/linear_model/_logistic.py:458: ConvergenceWarning: lbfgs failed to converge (status=
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.
```

```
Increase the number of iterations (max_iter) or scale the data as shown in:
    https://scikit-learn.org/stable/modules/preprocessing.html
Please also refer to the documentation for alternative solver options:
    https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression
  n_iter_i = _check_optimize_result(
/usr/local/lib/python3.10/dist-packages/sklearn/linear_model/_logistic.py:458: ConvergenceWarning: lbfgs failed to converge (status=
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.

Increase the number of iterations (max_iter) or scale the data as shown in:
    https://scikit-learn.org/stable/modules/preprocessing.html
Please also refer to the documentation for alternative solver options:
    https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression
  n_iter_i = _check_optimize_result(
/usr/local/lib/python3.10/dist-packages/sklearn/linear_model/_logistic.py:458: ConvergenceWarning: lbfgs failed to converge (status=
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.

Increase the number of iterations (max_iter) or scale the data as shown in:
    https://scikit-learn.org/stable/modules/preprocessing.html
Please also refer to the documentation for alternative solver options:
    https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression
  n_iter_i = _check_optimize_result(
/usr/local/lib/python3.10/dist-packages/sklearn/linear_model/_logistic.py:458: ConvergenceWarning: lbfgs failed to converge (status=
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.

Increase the number of iterations (max_iter) or scale the data as shown in:
    https://scikit-learn.org/stable/modules/preprocessing.html
Please also refer to the documentation for alternative solver options:
    https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression
  n_iter_i = _check_optimize_result(
/usr/local/lib/python3.10/dist-packages/sklearn/linear_model/_logistic.py:458: ConvergenceWarning: lbfgs failed to converge (status=
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.

Increase the number of iterations (max_iter) or scale the data as shown in:
    https://scikit-learn.org/stable/modules/preprocessing.html
Please also refer to the documentation for alternative solver options:
    https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression
  n_iter_i = _check_optimize_result(
/usr/local/lib/python3.10/dist-packages/sklearn/linear_model/_logistic.py:458: ConvergenceWarning: lbfgs failed to converge (status=
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.
```

## ▾ Model TPOT

```
 1 # Assuming you have the 'data1' and 'data2' DataFrames
 2 data1 = crypto_usd.copy()
 3 data2 = tweets.copy()
 4 # Merge the two DataFrames based on 'time' and 'date' columns
 5 merge = pd.merge(data1, data2, left_on='time', right_on='date')
 6
 7 # Drop the duplicate 'date' column
 8 merge.drop('date', axis=1, inplace=True)
 9
10 # Display the merged DataFrame
11 print(merge)
12
```

```
                    time     close      high       low      open  volumefrom  \
0     2023-02-25 21:00:00  22944.16  22960.69  22863.96  22921.71     1331.05
1     2023-02-25 21:00:00  22944.16  22960.69  22863.96  22921.71     1331.05
2     2023-02-25 21:00:00  22944.16  22960.69  22863.96  22921.71     1331.05
3     2023-02-25 21:00:00  22944.16  22960.69  22863.96  22921.71     1331.05
4     2023-02-25 21:00:00  22944.16  22960.69  22863.96  22921.71     1331.05
...                   ...       ...       ...       ...       ...         ...
7893  2023-03-04 23:00:00  22351.08  22352.28  22302.56  22311.46      476.12
7894  2023-03-04 23:00:00  22351.08  22352.28  22302.56  22311.46      476.12
7895  2023-03-04 23:00:00  22351.08  22352.28  22302.56  22311.46      476.12
7896  2023-03-04 23:00:00  22351.08  22352.28  22302.56  22311.46      476.12
7897  2023-03-04 23:00:00  22351.08  22352.28  22302.56  22311.46      476.12

         volumeto        Date      Time      volume  ...  user_verified  \
0     30505954.61  2023-02-25  21:00:00  30504623.56  ...          False
1     30505954.61  2023-02-25  21:00:00  30504623.56  ...          False
2     30505954.61  2023-02-25  21:00:00  30504623.56  ...          False
3     30505954.61  2023-02-25  21:00:00  30504623.56  ...          False
4     30505954.61  2023-02-25  21:00:00  30504623.56  ...          False
...           ...         ...       ...          ...  ...            ...
7893  10632637.83  2023-03-04  23:00:00  10632161.71  ...          False
7894  10632637.83  2023-03-04  23:00:00  10632161.71  ...          False
7895  10632637.83  2023-03-04  23:00:00  10632161.71  ...          False
```

```
7896  10632637.83  2023-03-04  23:00:00  10632161.71  ...         False
7897  10632637.83  2023-03-04  23:00:00  10632161.71  ...         False

                                              text  \
0     ethereum price updat eth 157128 usd bitcoin 00...
1              bitcoin 1month predict tuhgbqklxn
2     btcusdt 15m volum spike btc btc bitcoin ucl5iaaq4
3     lð□□□k take time think littlebit person load a...
4     ð□□□ð□□, sat 25 feb 2023 210035 gmt top 10 btc...
...                                             ...
7893  usd racist built colonist slaver paid btc bc e...
7894  everris rise everrisev3 everrevok defi crypto ...
7895  ð□□□ parti time ð□□□ ð□□□ 10000 x1 megapr ð□ª©...
7896  strategi 5010hl1h atr20d 92138 04 mar 2023 230...
7897  complet variou task hh8vl67nz5 claim slm token...

                                          hashtags               source  \
0     ['Ethereum', 'ETH', 'Bitcoin', 'BTC', 'altcoin...    Twitter Web App
1                                      ['Bitcoin']        predictCCbot
2                               ['BTC', 'Bitcoin']       JumpLineAlerts
3     ['GGA', 'cryptocurrency', 'Bitcoin', 'bnb', 'T...  Twitter for Android
4                                      ['bitcoin']              eht10c
...                                             ...                 ...
7893                                      ['BTC']  Twitter for Android
7894  ['EverRise', 'EverRiseV3', 'EverRevoke', 'DeFi...  EverRiseTwitterBot1
7895  ['btc', 'eth', 'xrp', 'doge', 'shiba', 'lto', ...    Twitter Web App
7896                              ['BTC', 'BitMEX']       system'cRe5520'
7897  ['SLMGames', 'SLM', 'Web3', 'BTC', 'ETH', 'BSC...           TweetDeck

      is_retweet  compound         score  sentiment_level  polarity  \
0            0.0    0.0000  0.000000e+00          Neutral  0.000000
1            0.0    0.0000  0.000000e+00          Neutral  0.000000
2            0.0    0.0000  0.000000e+00          Neutral  0.000000
3            0.0   -0.3089 -9.666133e+05         Negative -0.041667
4            0.0    0.2023  7.485100e+00         Positive  0.500000
```
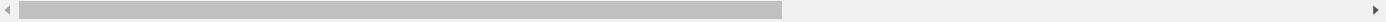
```
1 merge.head()
```

| | time | close | high | low | open | volumefrom | volumeto | Date | Time | volume | ... | user_verified | text |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 2023-02-25 21:00:00 | 22944.16 | 22960.69 | 22863.96 | 22921.71 | 1331.05 | 30505954.61 | 2023-02-25 | 21:00:00 | 30504623.56 | ... | False | ethereum price updat eth 157128 usd bitcoin 00... |
| 1 | 2023-02-25 21:00:00 | 22944.16 | 22960.69 | 22863.96 | 22921.71 | 1331.05 | 30505954.61 | 2023-02-25 | 21:00:00 | 30504623.56 | ... | False | bitcoin 1month predict tuhgbqklxn |
| 2 | 2023-02-25 21:00:00 | 22944.16 | 22960.69 | 22863.96 | 22921.71 | 1331.05 | 30505954.61 | 2023-02-25 | 21:00:00 | 30504623.56 | ... | False | btcusdt 15m volum spike btc btc bitcoin ucl5iaaq4 |
| 3 | 2023-02-25 21:00:00 | 22944.16 | 22960.69 | 22863.96 | 22921.71 | 1331.05 | 30505954.61 | 2023-02-25 | 21:00:00 | 30504623.56 | ... | False | lð□□□k take time think littlebit person load a... |
| 4 | 2023-02-25 21:00:00 | 22944.16 | 22960.69 | 22863.96 | 22921.71 | 1331.05 | 30505954.61 | 2023-02-25 | 21:00:00 | 30504623.56 | ... | False | ð□□□ð□□, sat 25 feb 2023 210035 gmt top 10 btc... |

5 rows × 29 columns

```
1 merge.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 7898 entries, 0 to 7897
Data columns (total 29 columns):
 #   Column          Non-Null Count  Dtype
```

```
 ---  ------           --------------  -----
  0   time             7898 non-null   object
  1   close            7898 non-null   float64
  2   high             7898 non-null   float64
  3   low              7898 non-null   float64
  4   open             7898 non-null   float64
  5   volumefrom       7898 non-null   float64
  6   volumeto         7898 non-null   float64
  7   Date             7898 non-null   object
  8   Time             7898 non-null   object
  9   volume           7898 non-null   float64
 10   marketcap        7898 non-null   float64
 11   price_delta      7898 non-null   float64
 12   user_name        7898 non-null   object
 13   user_location    3898 non-null   object
 14   user_description 7620 non-null   object
 15   user_created     7898 non-null   object
 16   user_followers   7898 non-null   float64
 17   user_friends     7898 non-null   float64
 18   user_favourites  7898 non-null   float64
 19   user_verified    7898 non-null   bool
 20   text             7898 non-null   object
 21   hashtags         7891 non-null   object
 22   source           7891 non-null   object
 23   is_retweet       7891 non-null   float64
 24   compound         7898 non-null   float64
 25   score            7898 non-null   float64
 26   sentiment_level  7898 non-null   object
 27   polarity         7898 non-null   float64
 28   subjectivity     7898 non-null   float64
dtypes: bool(1), float64(17), object(11)
memory usage: 1.8+ MB
```

```
1 label_counts = tweets['sentiment_level'].value_counts()
2 print(label_counts)
```

```
Neutral           93169
Positive          35921
Extreme Positive  17343
Negative          15903
Extreme Negative   5316
Name: sentiment_level, dtype: int64
```
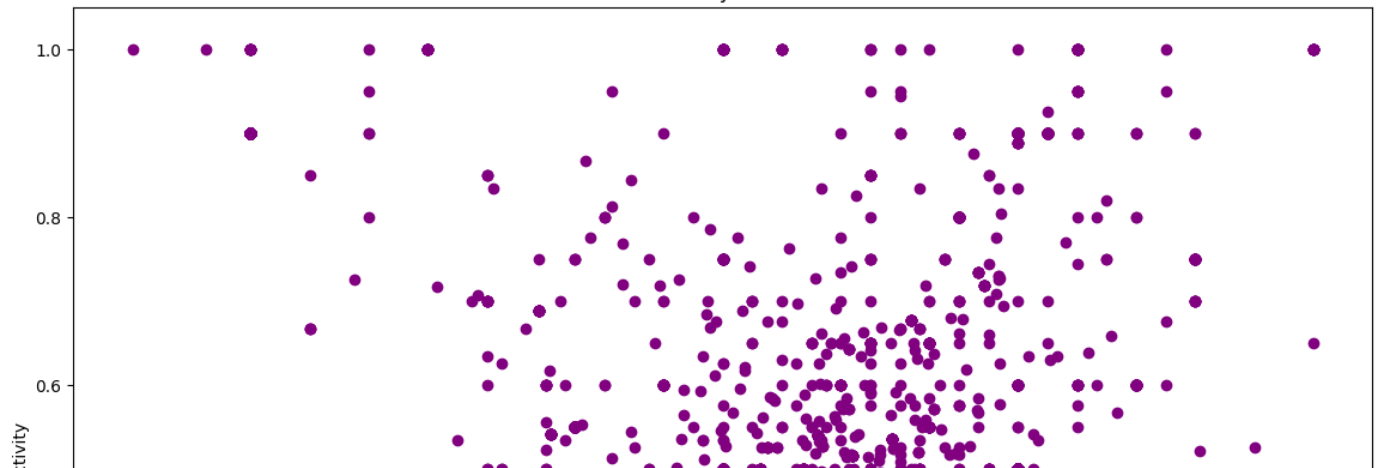
```
 1 import matplotlib.pyplot as plt
 2 # scatter plot to show the subjectivity and the polarity
 3 plt.figure(figsize=(14,10))
 4
 5 for i in range(merge.shape[0]):
 6     plt.scatter(merge["polarity"].iloc[[i]].values[0], merge["subjectivity"].iloc[[i]].values[0], color="Purple")
 7
 8 plt.title("Sentiment Analysis Scatter Plot")
 9 plt.xlabel('polarity')
10 plt.ylabel('subjectivity')
11 plt.show()
```

Sentiment Analysis Scatter Plot



```
1 #Creating Target Column
```

```
1 price_indicator = [merge.close[0] - merge['open'][0]]
2 for i in range(99):
3     price_indicator.append(merge.close[i+1] - merge.close[i])
4 #price_indicator
```

```
1 merge['price_indicator'] = 0
2 for i in range(len(price_indicator)):
3     merge['price_indicator'][i] = price_indicator[i]
4
5 merge.head()
```

```
<ipython-input-18-8f90b0759c32>:3: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-
  merge['price_indicator'][i] = price_indicator[i]
```

| | time | close | high | low | open | volumefrom | volumeto | Date | Time | volume | ... | text | hashtags |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 2023-02-25 21:00:00 | 22944.16 | 22960.69 | 22863.96 | 22921.71 | 1331.05 | 30505954.61 | 2023-02-25 | 21:00:00 | 30504623.56 | ... | ethereum price updat eth 157128 usd bitcoin 00... | ['Ethereum', 'ETH', 'Bitcoin', 'BTC', 'altcoin... |
| 1 | 2023-02-25 21:00:00 | 22944.16 | 22960.69 | 22863.96 | 22921.71 | 1331.05 | 30505954.61 | 2023-02-25 | 21:00:00 | 30504623.56 | ... | bitcoin 1month predict tuhgbqklxn | ['Bitcoin'] |
| 2 | 2023-02-25 21:00:00 | 22944.16 | 22960.69 | 22863.96 | 22921.71 | 1331.05 | 30505954.61 | 2023-02-25 | 21:00:00 | 30504623.56 | ... | btcusdt 15m volum spike btc btc bitcoin ucl5iaaq4 | ['BTC', 'Bitcoin'] |
| 3 | 2023-02-25 21:00:00 | 22944.16 | 22960.69 | 22863.96 | 22921.71 | 1331.05 | 30505954.61 | 2023-02-25 | 21:00:00 | 30504623.56 | ... | lð␣␣␣k take time think littlebit person load a... | ['GGA', 'cryptocurrency', 'Bitcoin', 'bnb', 'T... |
| 4 | 2023-02-25 21:00:00 | 22944.16 | 22960.69 | 22863.96 | 22921.71 | 1331.05 | 30505954.61 | 2023-02-25 | 21:00:00 | 30504623.56 | ... | ð␣␣␣ð␣␣¸ sat 25 feb 2023 210035 gmt top 10 btc... | ['bitcoin'] |

5 rows × 30 columns

```
1 merge['target'] = 0
2 for i in range(100):
3     if merge.price_indicator[i] > 0:
4         merge['target'][i] = 1
5
6 # 0 - price down
7 # 1 - price up
8
9 merge.head()
```

```
<ipython-input-19-e0c87f2219f9>:4: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame
```

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-
  merge['target'][i] = 1

| | time | close | high | low | open | volumefrom | volumeto | Date | Time | volume | ... | hashtags | source |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 2023-02-25 21:00:00 | 22944.16 | 22960.69 | 22863.96 | 22921.71 | 1331.05 | 30505954.61 | 2023-02-25 | 21:00:00 | 30504623.56 | ... | ['Ethereum', 'ETH', 'Bitcoin', 'BTC', 'altcoin... | Twitter Web App |
| 1 | 2023-02-25 21:00:00 | 22944.16 | 22960.69 | 22863.96 | 22921.71 | 1331.05 | 30505954.61 | 2023-02-25 | 21:00:00 | 30504623.56 | ... | ['Bitcoin'] | predictCCbo |
| 2 | 2023-02-25 21:00:00 | 22944.16 | 22960.69 | 22863.96 | 22921.71 | 1331.05 | 30505954.61 | 2023-02-25 | 21:00:00 | 30504623.56 | ... | ['BTC', 'Bitcoin'] | JumpLineAlerts |
| 3 | 2023-02-25 21:00:00 | 22944.16 | 22960.69 | 22863.96 | 22921.71 | 1331.05 | 30505954.61 | 2023-02-25 | 21:00:00 | 30504623.56 | ... | ['GGA', 'cryptocurrency', 'Bitcoin', 'bnb', 'T... | Twitter for Android |
| 4 | 2023-02-25 21:00:00 | 22944.16 | 22960.69 | 22863.96 | 22921.71 | 1331.05 | 30505954.61 | 2023-02-25 | 21:00:00 | 30504623.56 | ... | ['bitcoin'] | eht10 |

5 rows × 31 columns

```
1 keep_columns = ['open', 'high', 'low', 'close', 'volume','polarity','subjectivity','compound','score','price_indicator','target']
2 df = merge[keep_columns]
3 df.head()
```

| | open | high | low | close | volume | polarity | subjectivity | compound | score | price_indicator | target |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 22921.71 | 22960.69 | 22863.96 | 22944.16 | 30504623.56 | 0.000000 | 0.250000 | 0.0000 | 0.0000 | 22.45 | 1 |
| 1 | 22921.71 | 22960.69 | 22863.96 | 22944.16 | 30504623.56 | 0.000000 | 0.000000 | 0.0000 | 0.0000 | 0.00 | 0 |
| 2 | 22921.71 | 22960.69 | 22863.96 | 22944.16 | 30504623.56 | 0.000000 | 0.000000 | 0.0000 | 0.0000 | 0.00 | 0 |
| 3 | 22921.71 | 22960.69 | 22863.96 | 22944.16 | 30504623.56 | -0.041667 | 0.458333 | -0.3089 | -966613.2779 | 0.00 | 0 |
| 4 | 22921.71 | 22960.69 | 22863.96 | 22944.16 | 30504623.56 | 0.500000 | 0.500000 | 0.2023 | 7.4851 | 0.00 | 0 |

```
1 #Model Building
```

```
1 import numpy as np
2 #Create the feature data set
3 X = df
4 X = np.array(X.drop(['target'],1))
5 #Create the target data set
6 y = np.array(df['target'])
```

```
<ipython-input-22-63d9de6a3c5f>:4: FutureWarning: In a future version of pandas all arguments of DataFrame.drop except for the argumen
  X = np.array(X.drop(['target'],1))
```

```
1 from sklearn.model_selection import train_test_split
2 #Split the data into 80% training and 20% testing data sets
3 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state = 0)
```

```
1 !pip install tpot
2
```

```
Collecting tpot
  Downloading TPOT-0.12.0-py3-none-any.whl (87 kB)
  ━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 87.4/87.4 kB 5.8 MB/s eta 0:00:00
Requirement already satisfied: numpy>=1.16.3 in /usr/local/lib/python3.10/dist-packages (from tpot) (1.22.4)
Requirement already satisfied: scipy>=1.3.1 in /usr/local/lib/python3.10/dist-packages (from tpot) (1.10.1)
Requirement already satisfied: scikit-learn>=0.22.0 in /usr/local/lib/python3.10/dist-packages (from tpot) (1.2.2)
Collecting deap>=1.2 (from tpot)
  Downloading deap-1.3.3-cp310-cp310-manylinux_2_5_x86_64.manylinux1_x86_64.manylinux_2_17_x86_64.manylinux2014_x86_64.whl (139 kB)
  ━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 139.9/139.9 kB 14.0 MB/s eta 0:00:00
Collecting update-checker>=0.16 (from tpot)
  Downloading update_checker-0.18.0-py3-none-any.whl (7.0 kB)
Requirement already satisfied: tqdm>=4.36.1 in /usr/local/lib/python3.10/dist-packages (from tpot) (4.65.0)
Collecting stopit>=1.1.1 (from tpot)
  Downloading stopit-1.1.2.tar.gz (18 kB)
  Preparing metadata (setup.py) ... done
Requirement already satisfied: pandas>=0.24.2 in /usr/local/lib/python3.10/dist-packages (from tpot) (1.5.3)
Requirement already satisfied: joblib>=0.13.2 in /usr/local/lib/python3.10/dist-packages (from tpot) (1.2.0)
Requirement already satisfied: xgboost>=1.1.0 in /usr/local/lib/python3.10/dist-packages (from tpot) (1.7.6)
Requirement already satisfied: python-dateutil>=2.8.1 in /usr/local/lib/python3.10/dist-packages (from pandas>=0.24.2->tpot) (2.8.2)
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.10/dist-packages (from pandas>=0.24.2->tpot) (2022.7.1)
Requirement already satisfied: threadpoolctl>=2.0.0 in /usr/local/lib/python3.10/dist-packages (from scikit-learn>=0.22.0->tpot) (3.1.
Requirement already satisfied: requests>=2.3.0 in /usr/local/lib/python3.10/dist-packages (from update-checker>=0.16->tpot) (2.27.1)
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.10/dist-packages (from python-dateutil>=2.8.1->pandas>=0.24.2->tpot)
Requirement already satisfied: urllib3<1.27,>=1.21.1 in /usr/local/lib/python3.10/dist-packages (from requests>=2.3.0->update-checker>
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.10/dist-packages (from requests>=2.3.0->update-checker>=0.
Requirement already satisfied: charset-normalizer~=2.0.0 in /usr/local/lib/python3.10/dist-packages (from requests>=2.3.0->update-chec
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.10/dist-packages (from requests>=2.3.0->update-checker>=0.16->tp
Building wheels for collected packages: stopit
  Building wheel for stopit (setup.py) ... done
  Created wheel for stopit: filename=stopit-1.1.2-py3-none-any.whl size=11938 sha256=a56fda5b968cc0cd8d28799e3e03a41bd1d28bf77cd34f246
  Stored in directory: /root/.cache/pip/wheels/af/f9/87/bf5b3d565c2a007b4dae9d8142dccc85a9f164e517062dd519
Successfully built stopit
Installing collected packages: stopit, deap, update-checker, tpot
Successfully installed deap-1.3.3 stopit-1.1.2 tpot-0.12.0 update-checker-0.18.0
```

```
1 from tpot import TPOTClassifier
2 from sklearn.metrics import confusion_matrix,accuracy_score,roc_auc_score
```

```
1 from sklearn.metrics import roc_auc_score
2 from tpot import TPOTClassifier
3 import numpy as np
4
```

```
1 from sklearn.metrics import roc_auc_score
2 from tpot import TPOTClassifier
3 import numpy as np
4
5 # Instantiate TPOTClassifier
6 tpot = TPOTClassifier(
7     generations=5,
8     population_size=20,
9     verbosity=2,
10     scoring='roc_auc',
11     random_state=42,
12     disable_update_check=True,
13     config_dict='TPOT light'
14 )
15
16 # Convert X_train and y_train to NumPy arrays
17 X_train = np.array(X_train)
18 y_train = np.array(y_train)
19
20 # Ensure that there are at least two classes in y_train
21 if len(np.unique(y_train)) < 2:
22     raise ValueError("At least two classes are required in y_train for ROC AUC score calculation.")
23
24 try:
25     # Fit TPOTClassifier
26     tpot.fit(X_train, y_train)
27
28     # AUC score for tpot model
29     X_test = np.array(X_test)  # Assuming you have X_test data
30     y_test = np.array(y_test)  # Assuming you have y_test data
```

```
31
32      # Ensure that there are at least two classes in y_test
33      if len(np.unique(y_test)) < 2:
34          raise ValueError("At least two classes are required in y_test for ROC AUC score calculation.")
35
36      tpot_auc_score = roc_auc_score(y_test, tpot.predict_proba(X_test)[:, 1])
37      print(f'\nAUC score: {tpot_auc_score:.4f}')
38
39      # Print best pipeline steps
40      print('\nBest pipeline steps:')
41      for idx, (name, transform) in enumerate(tpot.fitted_pipeline_.steps, start=1):
42          print(f'{idx}. {transform}')
43
44 except ValueError as e:
45      print("Error:", str(e))
46
```

Optimization Progress: 33%                                   40/120 [00:11<00:29, 2.76pipeline/s]

```
/usr/local/lib/python3.10/dist-packages/sklearn/model_selection/_split.py:700: UserWarning: The least populated class in y has only 1
    warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/model_selection/_split.py:700: UserWarning: The least populated class in y has only 1
    warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/model_selection/_split.py:700: UserWarning: The least populated class in y has only 1
    warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/model_selection/_split.py:700: UserWarning: The least populated class in y has only 1
    warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/model_selection/_split.py:700: UserWarning: The least populated class in y has only 1
    warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/model_selection/_split.py:700: UserWarning: The least populated class in y has only 1
    warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/model_selection/_split.py:700: UserWarning: The least populated class in y has only 1
    warnings.warn(
Error: Only one class present in y_true. ROC AUC score is not defined in that case.
/usr/local/lib/python3.10/dist-packages/sklearn/model_selection/_split.py:700: UserWarning: The least populated class in y has only 1
    warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/model_selection/_split.py:700: UserWarning: The least populated class in y has only 1
    warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/model_selection/_split.py:700: UserWarning: The least populated class in y has only 1
    warnings.warn(
```

```
 1 # Instantiate TPOTClassifier
 2 tpot = TPOTClassifier(
 3      generations=5, #number of iterations to run ; pipeline optimisation process ; by default value is 100
 4      population_size=20, #number of individuals to retrain in the genetic programing popluation in every generation, by default value is 1
 5      verbosity=2, #it will state how much info TPOT will communicate while it is running
 6      scoring='roc_auc', #use to evaluate the quality of given pipeline
 7      random_state=42,
 8      disable_update_check=True,
 9      config_dict='TPOT light'
10 )
11 tpot.fit(X_train, y_train)
12
13 # AUC score for tpot model
14 tpot_auc_score = roc_auc_score(y_test, tpot.predict_proba(X_test)[:, 1])
15 print(f'\nAUC score: {tpot_auc_score:.4f}')
16
17 # Print best pipeline steps
18 print('\nBest pipeline steps:', end='\n')
19 for idx, (name, transform) in enumerate(tpot.fitted_pipeline_.steps, start=1):
20      # Print idx and transform
21      print(f'{idx}. {transform}')
```

Optimization Progress: 33%                                         40/120 [00:06<00:21, 3.71pipeline/s]

```
/usr/local/lib/python3.10/dist-packages/sklearn/model_selection/_split.py:700: UserWarning: The least populated class in y has only 1
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/model_selection/_split.py:700: UserWarning: The least populated class in y has only 1
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/model_selection/_split.py:700: UserWarning: The least populated class in y has only 1
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/model_selection/_split.py:700: UserWarning: The least populated class in y has only 1
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/model_selection/_split.py:700: UserWarning: The least populated class in y has only 1
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/model_selection/_split.py:700: UserWarning: The least populated class in y has only 1
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/model_selection/_split.py:700: UserWarning: The least populated class in y has only 1
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/model_selection/_split.py:700: UserWarning: The least populated class in y has only 1
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/model_selection/_split.py:700: UserWarning: The least populated class in y has only 1
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/model_selection/_split.py:700: UserWarning: The least populated class in y has only 1
  warnings.warn(
--------------------------------------------------------------------
IndexError                                Traceback (most recent call last)
/usr/local/lib/python3.10/dist-packages/tpot/base.py in fit(self, features, target, sample_weight, groups)
    816               warnings.simplefilter("ignore")
--> 817               self._pop, _ = eaMuPlusLambda(
    818                   population=self._pop,
```

≎ 26 frames

```
IndexError: tuple index out of range

During handling of the above exception, another exception occurred:

ValueError                                Traceback (most recent call last)
/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_ranking.py in _binary_roc_auc_score(y_true, y_score, sample_weight, max_fpr)
```

```
1 tpot.fitted_pipeline_
```

```
--> 339          raise ValueError(
```

Model 1: Decision tree classifier

```python
1 from sklearn.tree import DecisionTreeClassifier
2
3 clf = DecisionTreeClassifier(criterion='entropy', max_depth=8,
4                                      min_samples_leaf=10,
5                                      min_samples_split=6,
6                                      random_state=42)
7 clf.fit(X_train,y_train)
```

```
▼                         DecisionTreeClassifier
DecisionTreeClassifier(criterion='entropy', max_depth=8, min_samples_leaf=10,
                        min_samples_split=6, random_state=42)
```

```python
1 y_predicted = clf.predict(X_test)
```

```python
1 y_predicted
```

```
array([0, 0, 0, ..., 0, 0, 0])
```

```python
1
2 print( classification_report(y_test, y_predicted) )
```

```
              precision    recall  f1-score   support

           0       1.00      1.00      1.00      1580

    accuracy                           1.00      1580
   macro avg       1.00      1.00      1.00      1580
weighted avg       1.00      1.00      1.00      1580
```

```python
1 accuracy_score(y_test,y_predicted)*100
```

```
100.0
```

```
1 #Creating Pipeline to see which model has more accuracy
```

```
1 from sklearn.preprocessing import StandardScaler
2 from sklearn.decomposition import PCA
3 from sklearn.pipeline import Pipeline
4 from sklearn.linear_model import LogisticRegression
5 from sklearn.tree import DecisionTreeClassifier
6 from sklearn.ensemble import RandomForestClassifier
```

```
1 pipeline_lr = Pipeline([('scaler1',StandardScaler()),
2                         ('pca1',PCA(n_components=2)),
3                         ('lr_classifier',LogisticRegression(random_state=0))])
```

```
1 pipeline_dt = Pipeline([('scaler2',StandardScaler()),
2                         ('pca2',PCA(n_components=2)),
3                         ('dt_classifier',DecisionTreeClassifier())])
```

```
1 pipeline_randomforest = Pipeline([('scaler3',StandardScaler()),
2                         ('pca3',PCA(n_components=2)),
3                         ('rf_classifier',RandomForestClassifier())])
```

```
1 pipeline = [pipeline_lr,pipeline_dt,pipeline_randomforest]
```

```
1 best_accuracy=0.0
2 best_classifier=0
3 best_pipeline=""
```

```
1 pipe_dict = {0:'Logistic Regression', 1:'Decision Tree', 2:'RandomForest'}
2
3 for pipe in pipeline:
4     pipe.fit(X_train,y_train)
5
```

```
1 for i,model in enumerate(pipeline):
2     print("{}Test Accuracy: {}".format(pipe_dict[i],model.score(X_test,y_test)))
```

```
Logistic RegressionTest Accuracy: 1.0
Decision TreeTest Accuracy: 0.9987341772151899
RandomForestTest Accuracy: 1.0
```

✓ 21s    completed at 5:12 PM                                                                    ● ✕