‣ 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 | Date | Time |
|---|------|-------|------|-----|------|-----------|----------|------|------|
| 0 | 2023-02-19 13:00:00 | 24682.03 | 24715.82 | 24682.03 | 24707.39 | 903.97 | 22335943.28 | 2023-02-19 | 13:00:00 | 22 |

```
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_followers | user_friends | us |
|---|-----------|---------------|------------------|--------------|----------------|--------------|----|
| 0 | Irk | Vancouver, WA | Irk started investing in the stock market in 1... | 2018-08-11 03:17:00 | 116.0 | 8.0 | |
| 1 | Xiang Zhang | NaN | Professional Software Engineer ð□□»ð□□□Crypto ... | 2011-01-11 01:37:00 | 42.0 | 22.0 | |

```
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. C
```

```
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
```

|          | precision | recall | f1-score | support |
|----------|-----------|--------|----------|---------|
| 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 the
```

## 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: 0.859632616414193
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 d
```

## 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.svm import LinearSVC
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 SVM classifier
19     classifier = LinearSVC()
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: 0.870253164556962
Precision: 0.8670628029958947
Recall: 0.870253164556962
F1-Score: 0.8670875346312075
Predicted sentiment: ['Neutral']
                   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
```

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

### 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 analysis
```

## ▾ 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
```

```
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=1
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=1
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(
Model: Logistic Regression
Best Parameters: {'C': 10.0}
Cross-Validation Accuracy: 0.8732179009190496
Accuracy: 0.8816455696202532

Model: Gradient Boosting
Best Parameters: {'n_estimators': 300}
Cross-Validation Accuracy: 0.8779677430670395
Accuracy: 0.8822784810126583
```

## ‣ Model TPOT

[ ] ↳ *34 cells hidden*

✓  0s    completed at 4:20 PM                                                    ● ✕