

▼ Datasets

```

1 #Price

1 import pandas as pd
2
3 # URL to the raw CSV file
4 url = 'https://raw.githubusercontent.com/Amarpreet3/CIND-820-CAPSTONE/main/Sentimental%20Analysis/BitcoinPricePreprocessed.csv'
5
6 # Read the CSV file from the URL
7 crypto_usd = pd.read_csv(url)
8
9 # Display the first few rows of the data
10 print(crypto_usd.head())
11
12

```

		time	close	high	low	open	volume	from \
0	2023-02-19	13:00:00	24682.03	24715.82	24682.03	24707.39	903.97	
1	2023-02-19	14:00:00	24765.79	24792.85	24679.21	24682.03	1220.29	
2	2023-02-19	15:00:00	24928.21	25022.49	24751.96	24765.79	5074.50	
3	2023-02-19	16:00:00	24786.44	25175.28	24704.53	24928.21	7094.72	
4	2023-02-19	17:00:00	24364.95	24806.64	24346.17	24786.44	6896.84	

	volume	to	Date	Time	volume	marketcap	price_delta
0	2.233594e+07		2023-02-19	13:00:00	2.233504e+07	5.512964e+11	NaN
1	3.020300e+07		2023-02-19	14:00:00	3.020178e+07	7.480012e+11	83.76
2	1.263085e+08		2023-02-19	15:00:00	1.263034e+08	3.148644e+12	162.42
3	1.770671e+08		2023-02-19	16:00:00	1.770600e+08	4.388863e+12	-141.77
4	1.693379e+08		2023-02-19	17:00:00	1.693310e+08	4.125910e+12	-421.49

```

1 import pandas as pd
2
3 file_urls = [
4     'https://github.com/Amarpreet3/CIND-820-CAPSTONE/raw/main/Sentimental%20Analysis/BitcoinTweetsPreprocessed_1.csv',
5     'https://github.com/Amarpreet3/CIND-820-CAPSTONE/raw/main/Sentimental%20Analysis/BitcoinTweetsPreprocessed_2.csv',
6     'https://github.com/Amarpreet3/CIND-820-CAPSTONE/raw/main/Sentimental%20Analysis/BitcoinTweetsPreprocessed_3.csv',
7     'https://github.com/Amarpreet3/CIND-820-CAPSTONE/raw/main/Sentimental%20Analysis/BitcoinTweetsPreprocessed_4.csv',
8     'https://github.com/Amarpreet3/CIND-820-CAPSTONE/raw/main/Sentimental%20Analysis/BitcoinTweetsPreprocessed_5.csv',
9     'https://github.com/Amarpreet3/CIND-820-CAPSTONE/raw/main/Sentimental%20Analysis/BitcoinTweetsPreprocessed_6.csv'
10 ]
11
12 dfs = []
13
14 for url in file_urls:
15     # Read the CSV file
16     df = pd.read_csv(url)
17
18     # Append the DataFrame to the list
19     dfs.append(df)
20
21 # Combine all DataFrames into a single DataFrame
22 combined_df = pd.concat(dfs)
23
24 # Display the first few rows of the combined DataFrame
25 print(combined_df.head())
26

```

	user_name	user_location	\
0	Irk	Vancouver, WA	
1	Xiang Zhang	NaN	
2	Rhizoo	NaN	
3	Hari Marquez	Las Vegas, NV	
4	Bitcoin Candle Bot	Brazil	

	user_description	user_created	\
0	Irk started investing in the stock market in 1...	2018-08-11 03:17:00	
1	Professional Software Engineer 00000000Crypto ...	2011-01-11 01:37:00	
2	researcher. local maxima dunning&kruger spec...	2019-04-03 18:09:00	
3	Don't trust, verify. #Bitcoin El Salvador ...	2014-01-17 23:04:00	
4	Robot that posts the closure of the bitcoin da...	2021-01-06 01:36:00	

	user_followers	user_friends	user_favourites	user_verified	\
0	116.0	8.0	4580.0	False	

```
1         42.0         22.0         5.0         False
2         778.0        627.0        3205.0        False
3         222.0        521.0        13052.0       False
4         40.0         4.0          1.0          False

        date                                     text \
0  2023-02-25 23:59:00 bitcoin btc rest crypto ye bitcoin cryptocurr ...
1  2023-02-25 23:59:00 retriev invest fund current ongo tidexcoin kic...
2  2023-02-25 23:59:00 bull save monthli thread today good shit bitco...
3  2023-02-25 23:59:00 el salvador shape futur bitcoin membv32cn
4  2023-02-25 23:59:00 candl day 25022023 close open 2319406 high 232...

        hashtags                                source \
0      ['Bitcoin', 'crypto', 'NeedsMoreCrash']  Twitter Web App
1  ['Tidexcoin', 'Kicurrency', 'LMY', 'GMK', 'SYR...  Twitter for iPhone
2      ['bitcoin']                               Twitter Web App
3      ['Bitcoin']                               Twitter Web App
4      ['Bitcoin', 'Candle', 'BearMarket']       Bitcoin Candle Bot

    is_retweet  compound          score sentiment_level  polarity  subjectivity
0           0.0   -0.4019 -2.154092e+05      Negative  0.000000    0.000000
1           0.0   0.0000  0.000000e+00       Neutral  0.000000    0.400000
2           0.0   0.3612  9.005682e+06       Positive  0.250000    0.700000
3           0.0   0.0000  0.000000e+00       Neutral  0.000000    0.000000
4           0.0  -0.2732 -2.240240e+01      Negative  0.053333    0.446667
```

```
1 tweets = combined_df.copy()
```

```
1 tweets.head()
```

	user_name	user_location	user_description	user_created	user_followers	user_friends	u:
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 ðŸ’Ž»ðŸ’ŽCrypto ...	2011-01-11 01:37:00	42.0	22.0	
2	Rhizoo	NaN	researcher. local maxima dunningâŸ’kruger spec...	2019-04-03 18:09:00	778.0	627.0	
3	Hari Marquez	Las Vegas, NV	DonâŸ’t trust, verify. #Bitcoin El Salvador ...	2014-01-17 23:04:00	222.0	521.0	
4	Bitcoin Candle Bot	Brazil	Robot that posts the closure of the bitcoin da...	2021-01-06 01:36:00	40.0	4.0	

```
1 print(tweets.columns)
```

```
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'],  
      dtype='object')
```

```
1 import pandas as pd
2
3
4 # Check the shape of the dataset
5 print("Shape of the dataset:", tweets.shape)
6
7 # Check the size of the dataset
```

```
8 print("Size of the dataset (number of elements):", tweets.size)
9
```

```
Shape of the dataset: (167652, 18)
Size of the dataset (number of elements): 3017736
```

```
1 import pandas as pd
2 import os
3
4
5 # Check the shape of the data
6 print("Shape of the data:", tweets.shape)
7
8
```

```
Shape of the data: (167652, 18)
```

```
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
```

▼ 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	u:
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 ðŸ’Ž»ðŸ’Ž»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 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

▼ 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

▼ Random Forest

```

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

▼ Logistic Regression:

```

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

▼ 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

▼ 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
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000

```

```

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

```

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...

```

```
/896 ['BIC', 'BITMEX'] system:ckes520
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	
...	
7893	0.0	-0.6124	-3.007276e+06	Extreme Negative	0.000000	
7894	0.0	0.0000	0.000000e+00	Neutral	0.000000	
7895	0.0	0.0000	0.000000e+00	Neutral	0.000000	
7896	0.0	0.2732	4.847934e+03	Positive	0.000000	
7897	0.0	0.8126	1.214187e+04	Extreme Positive	0.000000	

	subjectivity
0	0.250000
1	0.000000
2	0.000000
3	0.458333
4	0.500000
...	...
7893	0.000000
7894	0.000000
7895	0.000000
7896	0.000000
7897	0.000000

[7898 rows x 29 columns]

```
1 merge.head()
```

	time	close	high	low	open	volumefrom	volumeto	Date	Time
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 30
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 30
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 30
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 30
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 30

5 rows x 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)

```

```

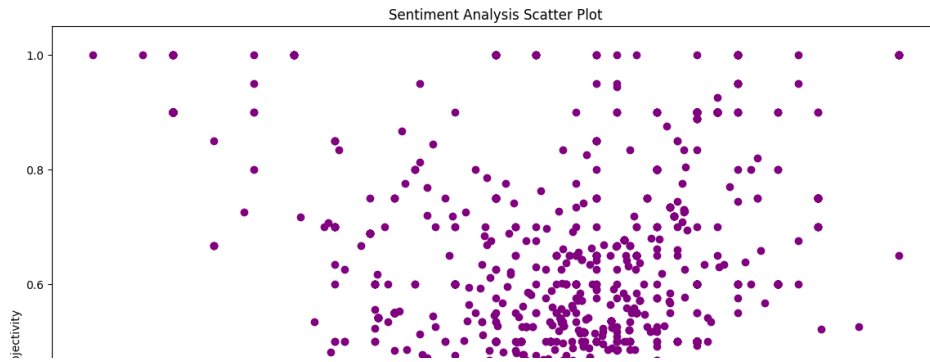
Neutral          93169
Positive         35921
Extreme Positive 17343
Negative         15903
Extreme Negative  5316
Name: sentiment_level, dtype: int64

```

```

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()

```



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-23-8f90b0759c32>: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
merge['price_indicator'][i] = price_indicator[i]

	time	close	high	low	open	volumefrom	volumeto	Date	Time
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 30
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 30
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 30
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 30
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 30

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-24-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

```
merge['target'][i] = 1
```

	time	close	high	low	open	volumefrom	volumeto	Date	Time
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 30
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 30
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 30
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 30
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 30

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
0	22921.71	22960.69	22863.96	22944.16	30504623.56	0.000000	0.250000	0.0000
1	22921.71	22960.69	22863.96	22944.16	30504623.56	0.000000	0.000000	0.0000
2	22921.71	22960.69	22863.96	22944.16	30504623.56	0.000000	0.000000	0.0000
3	22921.71	22960.69	22863.96	22944.16	30504623.56	-0.041667	0.458333	-0.3089 -9666
4	22921.71	22960.69	22863.96	22944.16	30504623.56	0.500000	0.500000	0.2023

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-42-63d9de6a3c5f>:4: FutureWarning: In a future version of pandas all arguments of DataFrame.drop except for the argument
X = np.array(X.drop(['target'],1))
```

```
1 #Split the data into 80% training and 20% testing data sets
2 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
```

```
Requirement already satisfied: tpot in /usr/local/lib/python3.10/dist-packages (0.12.0)
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)
Requirement already satisfied: deap>=1.2 in /usr/local/lib/python3.10/dist-packages (from tpot) (1.3.3)
Requirement already satisfied: update-checker>=0.16 in /usr/local/lib/python3.10/dist-packages (from tpot) (0.18.0)
Requirement already satisfied: tqdm>=4.36.1 in /usr/local/lib/python3.10/dist-packages (from tpot) (4.65.0)
Requirement already satisfied: stopit>=1.1.1 in /usr/local/lib/python3.10/dist-packages (from tpot) (1.1.2)
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) (3.1.0)
 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>=0.16->tpot)
 Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.10/dist-packages (from requests>=2.3.0->update-checker>=0.16->tpot)
 Requirement already satisfied: charset-normalizer~=2.0.0 in /usr/local/lib/python3.10/dist-packages (from requests>=2.3.0->update-checker>=0.16->tpot)
 Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.10/dist-packages (from requests>=2.3.0->update-checker>=0.16->tpot)

```
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 # Fit TPOTClassifier
21 tpot.fit(X_train, y_train)
22
23 # AUC score for tpot model
24 X_test = np.array(X_test) # Assuming you have X_test data
25 y_test = np.array(y_test) # Assuming you have y_test data
26
27 tpot_auc_score = roc_auc_score(y_test, tpot.predict_proba(X_test)[: , 1])
28 print(f'\nAUC score: {tpot_auc_score:.4f}')
29
30 # Print best pipeline steps
31 print('\nBest pipeline steps:')
32 for idx, (name, transform) in enumerate(tpot.fitted_pipeline_.steps, start=1):
33     print(f'{idx}. {transform}')
34
```

```

Optimization Progress: 33%                                40/120 [00:11<00:47, 1.69pipeline/s]
/usr/local/lib/python3.10/dist-packages/sklearn/model_selection/_split.py:700: UserWarning:
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/model_selection/_split.py:700: UserWarning:
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/model_selection/_split.py:700: UserWarning:
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/model_selection/_split.py:700: UserWarning:
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/model_selection/_split.py:700: UserWarning:
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/model_selection/_split.py:700: UserWarning:
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/model_selection/_split.py:700: UserWarning:
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/model_selection/_split.py:700: UserWarning:
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/model_selection/_split.py:700: UserWarning:
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/model_selection/_split.py:700: UserWarning:
  warnings.warn(
-----
IndexError                                Traceback (most recent call last)
/usr/local/lib/python3.10/dist-packages/tpot/base.py in fit(self, features, target,

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 programming population in every generation, by default value is 100
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}')

1 tpot.fitted_pipeline_

```

Model 1: Decision tree classifier

```

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)

1 y_predicted = clf.predict(X_test)

1 y_predicted

1
2 print( classification_report(y_test, y_predicted) )

1 accuracy_score(y_test,y_predicted)*100

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)))
```