

```

1 import pandas as pd
2 import re
3
4 from sklearn.model_selection import train_test_split
5 from sklearn.feature_extraction.text import TfidfVectorizer
6 from sklearn.linear_model import LogisticRegression
7 from sklearn.metrics import accuracy_score, classification_report

```

```

1 df = pd.read_csv("BTC_tweets_daily_example.csv")
2 print(df.head())

```

```

Unnamed: 0      Date \
0      0  Fri Mar 23 00:40:32 +0000 2018
1      1  Fri Mar 23 00:40:34 +0000 2018
2      2  Fri Mar 23 00:40:35 +0000 2018
3      3  Fri Mar 23 00:40:36 +0000 2018
4      4  Fri Mar 23 00:40:36 +0000 2018

      Tweet      Screen_name \
0  RT @ALXTOKEN: Paul Krugman, Nobel Luddite. I h...  myresumerocket
1  @lopp @Kevin_Pham @psycho_sage @naval But @Pr...  BitMacro
2  RT @tippereconomy: Another use case for #block...  hojachotopur
3      free coins https://t.co/DiuoePJdap  denies_distro
4  RT @payvxofficial: WE are happy to announce th...  aditzgraha

      Source \
0      []
1      [u'Bitcoin']
2  [u'blockchain', u'Tipper', u'TipperEconomy']
3      []
4      []

      Link      Sentiment \
0  <a href="http://twitter.com" rel="nofollow">Tw...  ['neutral']
1  <a href="http://twitter.com/download/android" ...  ['neutral']
2  <a href="http://twitter.com" rel="nofollow">Tw...  ['positive']
3  <a href="http://twitter.com" rel="nofollow">Tw...  ['positive']
4  <a href="http://twitter.com/download/android" ...  ['positive']

sent_score  New_Sentiment_Score  New_Sentiment_State
0      0.0      0.000000      0.0
1      0.0      0.000000      0.0
2      1.0      0.136364      1.0
3      1.0      0.400000      1.0
4      1.0      0.468182      1.0

```

```

1 def clean_text(text):
2     text = text.lower()
3     text = re.sub(r"http\S+", "", text)
4     text = re.sub(r"@w+", "", text)
5     text = re.sub(r"#w+", "", text)
6     text = re.sub(r"^[a-z\s]", "", text)
7     return text
8

```

```

1 # Step 4: Apply cleaning
2 df["Tweet"] = df["Tweet"].astype(str)
3 df["Tweet"] = df["Tweet"].apply(clean_text)

```

```

1 # Step 5: Split data (IMPORTANT - avoids leakage)
2 # Remove rows where 'New_Sentiment_State' is NaN
3 df_cleaned = df.dropna(subset=['New_Sentiment_State'])
4
5 X = df_cleaned["Tweet"]
6 y = df_cleaned["New_Sentiment_State"]
7
8 X_train, X_test, y_train, y_test = train_test_split(
9     X, y, test_size=0.2, random_state=42
10 )
11

```

```

1 # Step 6: TF-IDF Vectorization
2 tfidf = TfidfVectorizer(stop_words="english", max_features=5000)
3
4 X_train_tfidf = tfidf.fit_transform(X_train) # fit ONLY on training data
5 X_test_tfidf = tfidf.transform(X_test)
6

```

```

1 # Step 7: Train Logistic Regression model
2 model = LogisticRegression(max_iter=1000)

```

```

2 model = LogisticRegression(max_iter=1000,
3
4 model.fit(X_train_tfidf, y_train)
5

```

▼ LogisticRegression ⓘ ?

LogisticRegression(max\_iter=1000)

```

1 # Step 8: Predictions
2 y_pred = model.predict(X_test_tfidf)
3

```

```

1 # Step 9: Evaluation
2 print("Accuracy:", accuracy_score(y_test, y_pred))
3 print("\nClassification Report:\n")
4 print(classification_report(y_test, y_pred))
5

```

Accuracy: 0.9338314816635532

Classification Report:

	precision	recall	f1-score	support
-1.0	0.97	0.82	0.89	1298
0.0	0.90	0.97	0.93	4305
1.0	0.96	0.93	0.95	4568
accuracy			0.93	10171
macro avg	0.94	0.91	0.92	10171
weighted avg	0.94	0.93	0.93	10171

```

1 # Step 10: Predict sentiment of new BTC tweet
2 sample_tweet = ["Charan is good boy"]
3 sample_clean = [clean_text(sample_tweet[0])]
4 sample_tfidf = tfidf.transform(sample_clean)
5 prediction = model.predict(sample_tfidf)
6 probability = model.predict_proba(sample_tfidf)
7 print("Predicted Sentiment:", prediction[0])
8 print("Probabilities:", probability)

```

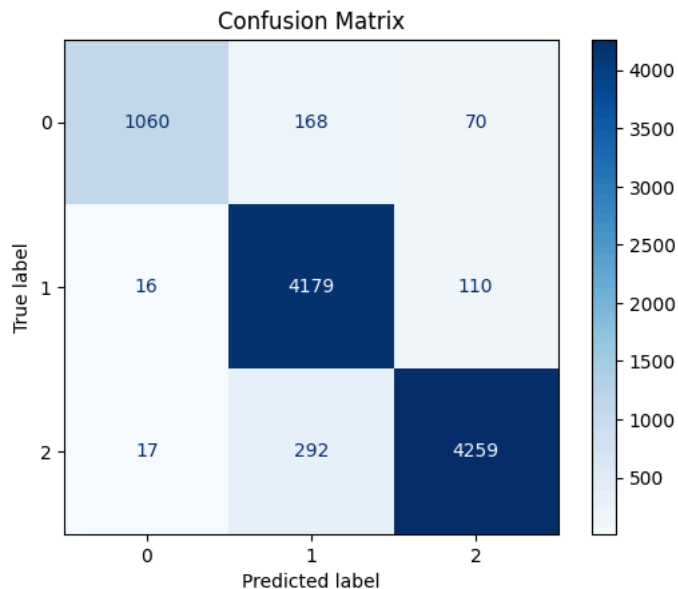
Predicted Sentiment: 1.0

Probabilities: [[0.00735989 0.02634762 0.96629249]]

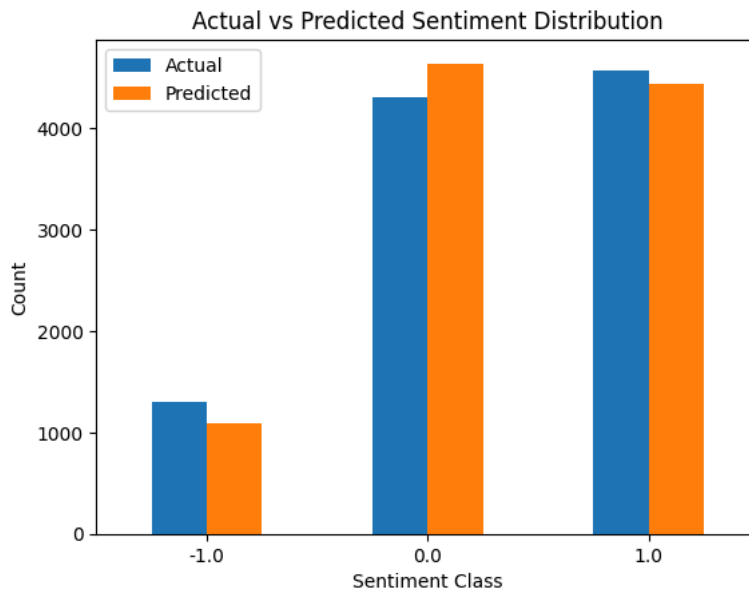
```

1 from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay
2 import matplotlib.pyplot as plt
3
4 cm = confusion_matrix(y_test, y_pred)
5
6 disp = ConfusionMatrixDisplay(confusion_matrix=cm)
7 disp.plot(cmap="Blues")
8 plt.title("Confusion Matrix")
9 plt.show()
10

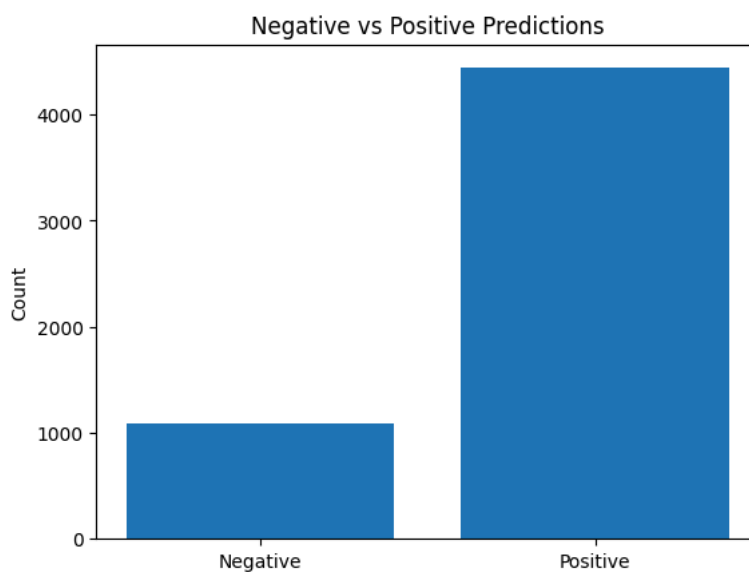
```



```
1 import pandas as pd
2
3 actual_counts = pd.Series(y_test).value_counts().sort_index()
4 pred_counts = pd.Series(y_pred).value_counts().sort_index()
5
6 df_plot = pd.DataFrame({
7     "Actual": actual_counts,
8     "Predicted": pred_counts
9 })
10
11 df_plot.plot(kind="bar")
12 plt.title("Actual vs Predicted Sentiment Distribution")
13 plt.xlabel("Sentiment Class")
14 plt.ylabel("Count")
15 plt.xticks(rotation=0)
16 plt.show()
17
```

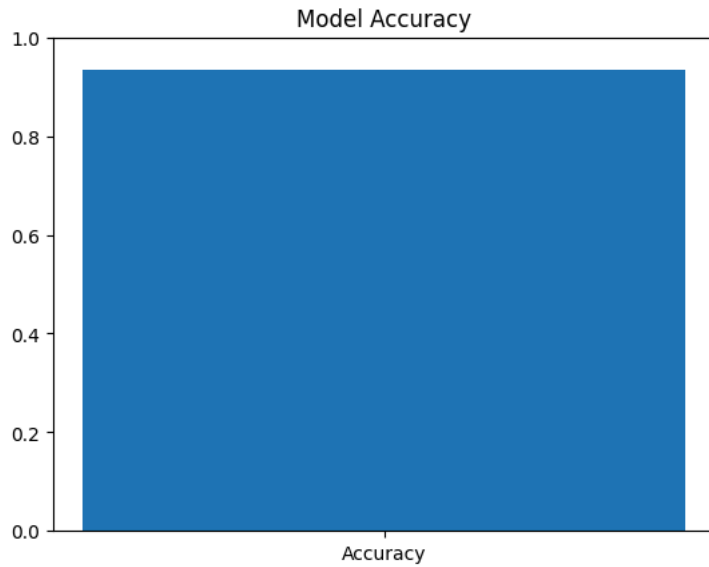


```
1 neg_pos = {
2     "Negative": pred_counts.get(-1, 0),
3     "Positive": pred_counts.get(1, 0)
4 }
5
6 plt.bar(neg_pos.keys(), neg_pos.values())
7 plt.title("Negative vs Positive Predictions")
8 plt.ylabel("Count")
9 plt.show()
10
```



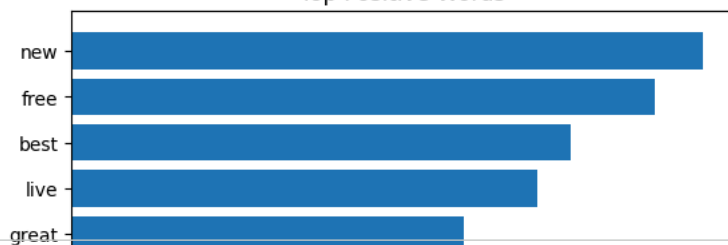
```
1 acc = accuracy_score(y_test, y_pred)
2
3 plt.bar(["Accuracy"], [acc])
```

```
4 plt.ylim(0, 1)
5 plt.title("Model Accuracy")
6 plt.show()
7
```



```
1 import numpy as np
2
3 feature_names = tfidf.get_feature_names_out()
4 coef = model.coef_
5
6 # POSITIVE class (highest weights)
7 top_positive = np.argsort(coef[2])[-10:]
8 plt.barh(feature_names[top_positive], coef[2][top_positive])
9 plt.title("Top Positive Words")
10 plt.show()
11
12 # NEGATIVE class (lowest weights)
13 top_negative = np.argsort(coef[0])[:10]
14 plt.barh(feature_names[top_negative], coef[0][top_negative])
15 plt.title("Top Negative Words")
16 plt.show()
17
```

Top Positive Words



```

1 # TF-IDF + Logistic Regression prediction
2
3 user_input = ["I just lost all my money in Bitcoin, great job crypto"]
4 user_clean = [clean_text(user_input[0])]
5
6 user_tfidf = tfidf.transform(user_clean)
7 tfidf_pred = model.predict(user_tfidf)
8
9 print("TF-IDF Prediction:", tfidf_pred[0])
10

```

TF-IDF Prediction: 1.0

Top Negative Words

```

1 # Assume TF-IDF vectorizer (tfidf) and logistic model (log_model) are already trained
2
3 sentence = "I just love how Bitcoin destroyed my savings today"
4 clean_sentence = clean_text(sentence)
5
6 X_vec = tfidf.transform([clean_sentence])
7 tfidf_pred = model.predict(X_vec)
8
9 print("TF-IDF + Logistic Prediction:", tfidf_pred[0])

```

TF-IDF + Logistic Prediction: 1.0

