```
1 # from google.colab import files
2 # uploaded = files.upload()
```

```
1 import pandas as pd
2 import numpy as np
```

```
1 rev_df = pd.read_csv('Amazon_Reviews.csv', engine='python')
2 rev_df
```

| | Reviewer Name | Profile Link | Country | Review Count | Review Date | Rating | Review Title | Review Text | Date of Experience |
|---|---|---|---|---|---|---|---|---|---|
| 0 | Eugene ath | /users/66e8185ff1598352d6b3701a | US | 1 review | 2024-09-16T13:44:26.000Z | Rated 1 out of 5 stars | A Store That Doesn't Want to Sell Anything | I registered on the website, tried to order a ... | September 16, 2024 |
| 1 | Daniel ohalloran | /users/5d75e460200c1f6a6373648c | GB | 9 reviews | 2024-09-16T18:26:46.000Z | Rated 1 out of 5 stars | Had multiple orders one turned up and… | Had multiple orders one turned up and driver h... | September 16, 2024 |
| 2 | p fisher | /users/546cfcf1000064000197b88f | GB | 90 reviews | 2024-09-16T21:47:39.000Z | Rated 1 out of 5 stars | I informed these reprobates | I informed these reprobates that I WOULD NOT B... | September 16, 2024 |
| 3 | Greg Dunn | /users/62c35cdbacc0ea0012ccaffa | AU | 5 reviews | 2024-09-17T07:15:49.000Z | Rated 1 out of 5 stars | Advertise one price then increase it on website | I have bought from Amazon before and no proble... | September 17, 2024 |
| 4 | Sheila Hannah | /users/5ddbe429478d88251550610e | GB | 8 reviews | 2024-09-16T18:37:17.000Z | Rated 1 out of 5 stars | If I could give a lower rate I would | If I could give a lower rate I would! I cancel... | September 16, 2024 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |

Next steps: ( Generate code with `rev_df` ) ( 🔘 View recommended plots ) ( New interactive sheet )

```
1 rev_df.info()
2 print()
3 print(f'Null : \n{rev_df.isnull().sum()}')
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 21214 entries, 0 to 21213
Data columns (total 9 columns):
 #   Column              Non-Null Count  Dtype
---  ------              --------------  -----
 0   Reviewer Name       21214 non-null  object
 1   Profile Link        21163 non-null  object
 2   Country             21054 non-null  object
 3   Review Count        21055 non-null  object
 4   Review Date         21055 non-null  object
 5   Rating              21055 non-null  object
 6   Review Title        21055 non-null  object
 7   Review Text         21055 non-null  object
 8   Date of Experience  20947 non-null  object
dtypes: object(9)
memory usage: 1.5+ MB

Null :
Reviewer Name         0
Profile Link         51
Country             160
Review Count        159
Review Date         159
Rating              159
Review Title        159
Review Text         159
Date of Experience  267
dtype: int64
```

```python
1 # correcting all the column names
2
3 rev_df.columns = rev_df.columns.str.lower().str.replace(' ', '_')
4 rev_df.columns
```

```
Index(['reviewer_name', 'profile_link', 'country', 'review_count',
       'review_date', 'rating', 'review_title', 'review_text',
       'date_of_experience'],
      dtype='object')
```

```python
1 # cleaning the profile links and removing duplicates
2 import uuid
3
4 uid = uuid.uuid4().hex[:24]
5 profile = f'/user/{uid}'
6 rev_df['profile_link'].fillna(profile,inplace=True)
7 rev_df.drop_duplicates(inplace=True)
8
```

```
<ipython-input-7-e2d94c52b8d0>:6: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained assi
  The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting

  For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method({col: value}, inplace=True)' or df[col] = df[col


    rev_df['profile_link'].fillna(profile,inplace=True)
```

```python
1 # usingmos common country to fill in the country names using mode
2 most_common_country = rev_df['country'].mode()[0]
3 rev_df['country'].fillna(most_common_country, inplace=True)
4 rev_df
```

```
<ipython-input-8-e7df5d195954>:3: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained assi
  The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting

  For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method({col: value}, inplace=True)' or df[col] = df[col


    rev_df['country'].fillna(most_common_country, inplace=True)
```

| | reviewer_name | profile_link | country | review_count | review_date | rating | review_title | review_text | dat |
|---|---|---|---|---|---|---|---|---|---|
| 0 | Eugene ath | /users/66e8185ff1598352d6b3701a | US | 1 review | 2024-09-16T13:44:26.000Z | Rated 1 out of 5 stars | A Store That Doesn't Want to Sell Anything | I registered on the website, tried to order a ... | Se |
| 1 | Daniel ohalloran | /users/5d75e460200c1f6a6373648c | GB | 9 reviews | 2024-09-16T18:26:46.000Z | Rated 1 out of 5 stars | Had multiple orders one turned up and… | Had multiple orders one turned up and driver h... | Se |
| 2 | p fisher | /users/546cfcf1000064000197b88f | GB | 90 reviews | 2024-09-16T21:47:39.000Z | Rated 1 out of 5 stars | I informed these reprobates | I informed these reprobates that I WOULD NOT B... | Se |
| 3 | Greg Dunn | /users/62c35cdbacc0ea0012ccaffa | AU | 5 reviews | 2024-09-17T07:15:49.000Z | Rated 1 out of 5 stars | Advertise one price then increase it on website | I have bought from Amazon before and no proble... | Se |
| 4 | Sheila Hannah | /users/5ddbe429478d88251550610e | GB | 8 reviews | 2024-09-16T18:37:17.000Z | Rated 1 out of 5 stars | If I could give a lower rate I would | If I could give a lower rate I would! I cancel... | Se |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 21209 | Anders T | /users/47bd4ffe0000640001001044 | DK | 1 review | 2009-03-22T13:14:12.000Z | Rated 5 out of 5 | Fast!! | I have had perfect order fulfillment, | Se |

Next steps:  ( Generate code with `rev_df` )  ( 💬 View recommended plots )  ( New interactive sheet )

```python
1 # from the review count column numerocal values are extraxcted
2
3 rev_df['review_count'] = rev_df['review_count'].str.extract(r'(\d+)')
4 rev_df['review_count'] = pd.to_numeric(rev_df['review_count'], errors='coerce')
5 mean_val = rev_df['review_count'].mean()
6 rev_df['review_count'].fillna(mean_val, inplace=True)
```

```
7 rev_df['review_count'] = rev_df['review_count'].astype(int)
8 rev_df
```

⟫ `<ipython-input-9-df4aa9abf6db>:6: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained assi`
`The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting`

`For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method({col: value}, inplace=True)' or df[col] = df[col]`

`rev_df['review_count'].fillna(mean_val, inplace=True)`

| | reviewer_name | profile_link | country | review_count | review_date | rating | review_title | review_text | dat |
|---|---|---|---|---|---|---|---|---|---|
| **0** | Eugene ath | /users/66e8185ff1598352d6b3701a | US | 1 | 2024-09-16T13:44:26.000Z | Rated 1 out of 5 stars | A Store That Doesn't Want to Sell Anything | I registered on the website, tried to order a ... | Se |
| **1** | Daniel ohalloran | /users/5d75e460200c1f6a6373648c | GB | 9 | 2024-09-16T18:26:46.000Z | Rated 1 out of 5 stars | Had multiple orders one turned up and… | Had multiple orders one turned up and driver h... | Se |
| **2** | p fisher | /users/546cfcf1000064000197b88f | GB | 90 | 2024-09-16T21:47:39.000Z | Rated 1 out of 5 stars | I informed these reprobates | I informed these reprobates that I WOULD NOT B... | Se |
| **3** | Greg Dunn | /users/62c35cdbacc0ea0012ccaffa | AU | 5 | 2024-09-17T07:15:49.000Z | Rated 1 out of 5 stars | Advertise one price then increase it on website | I have bought from Amazon before and no proble... | Se |
| **4** | Sheila Hannah | /users/5ddbe429478d88251550610e | GB | 8 | 2024-09-16T18:37:17.000Z | Rated 1 out of 5 stars | If I could give a lower rate I would | If I could give a lower rate I would! I cancel... | Se |
| **...** | ... | ... | ... | ... | ... | ... | ... | ... | |
| **21209** | Anders T | /users/47bd4ffe0000640001001044 | DK | 1 | 2009-03-22T13:14:12.000Z | Rated 5 out of 5 | Fast!! | I have had perfect order fulfillment, | |

Next steps: ( Generate code with `rev_df` ) ( ⬤ View recommended plots ) ( New interactive sheet )

```
1 # Transforming the review_date column as per datetime
2
3 rev_df['review_date'] = pd.to_datetime(rev_df['review_date'], errors='coerce')
4 # rev_df['review_date'] = rev_df['review_date'].dt.strftime('%Y-%m-%d')
5 most_common_date = rev_df['review_date'].mode()[0]
6 rev_df['review_date'] = rev_df['review_date'].fillna(most_common_date)
7
8 rev_df['date_of_experience'] = pd.to_datetime(rev_df['date_of_experience'])
9 # rev_df['date_of_experience'] = rev_df['date_of_experience'].dt.strftime('%Y-%m-%d')
10 most_common_date = rev_df['date_of_experience'].mode()[0]
11 rev_df['date_of_experience'] = rev_df['date_of_experience'].fillna(most_common_date)
12
13 rev_df
```

| | reviewer_name | profile_link | country | review_count | review_date | rating | review_title | review_text | date_o |
|---|---|---|---|---|---|---|---|---|---|
| 0 | Eugene ath | /users/66e8185ff1598352d6b3701a | US | 1 | 2024-09-16 13:44:26+00:00 | Rated 1 out of 5 stars | A Store That Doesn't Want to Sell Anything | I registered on the website, tried to order a ... | |
| 1 | Daniel ohalloran | /users/5d75e460200c1f6a6373648c | GB | 9 | 2024-09-16 18:26:46+00:00 | Rated 1 out of 5 stars | Had multiple orders one turned up and… | Had multiple orders one turned up and driver h... | |
| 2 | p fisher | /users/546cfcf1000064000197b88f | GB | 90 | 2024-09-16 21:47:39+00:00 | Rated 1 out of 5 stars | I informed these reprobates | I informed these reprobates that I WOULD NOT B... | |
| 3 | Greg Dunn | /users/62c35cdbacc0ea0012ccaffa | AU | 5 | 2024-09-17 07:15:49+00:00 | Rated 1 out of 5 stars | Advertise one price then increase it on website | I have bought from Amazon before and no proble... | |
| 4 | Sheila Hannah | /users/5ddbe429478d88251550610e | GB | 8 | 2024-09-16 18:37:17+00:00 | Rated 1 out of 5 stars | If I could give a lower rate I would | If I could give a lower rate I would! I cancel... | |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| | Andres T | /users/47b14ff... | DK | | 2009-03-22 | Rated 5 out | | I have had perfect order | |

Next steps:  [ Generate code with `rev_df` ]   [ 🔘 View recommended plots ]   [ New interactive sheet ]

```
1 # Extracting the numerical value from string
2
3 rev_df['rating'] = rev_df['rating'].str.extract(r'Rated (\d+) out of 5 stars')
4 rev_df['rating'] = pd.to_numeric(rev_df['rating'], errors='coerce')
5 mean_value = rev_df['rating'].mean()
6 rev_df['rating'] = rev_df['rating'].fillna(mean_value).round().astype(int)
```

```
1 rev_df['review_title'] = rev_df['review_title'].fillna('N/A')
2 rev_df['review_text'] = rev_df['review_text'].fillna('N/A')
3
4 rev_df['text_features'] = rev_df['review_text'] + ' ' + rev_df['review_title']
```

```
1 from textblob import TextBlob
2
3 rev_df['polarity'] = rev_df['review_text'].apply(lambda x: TextBlob(x).sentiment.polarity)
4 rev_df['sentiment'] = rev_df['polarity'].apply(lambda x: 'positive' if x > 0 else ('negative' if x < 0 else 'neutral'))
5 rev_df
```

| | reviewer_name | profile_link | country | review_count | review_date | rating | review_title | review_text | dat |
|---|---|---|---|---|---|---|---|---|---|
| **0** | Eugene ath | /users/66e8185ff1598352d6b3701a | US | 1 | 2024-09-16 13:44:26+00:00 | 1 | A Store That Doesn't Want to Sell Anything | I registered on the website, tried to order a ... | |
| **1** | Daniel ohalloran | /users/5d75e460200c1f6a6373648c | GB | 9 | 2024-09-16 18:26:46+00:00 | 1 | Had multiple orders one turned up and… | Had multiple orders one turned up and driver h... | |
| **2** | p fisher | /users/546cfcf1000064000197b88f | GB | 90 | 2024-09-16 21:47:39+00:00 | 1 | I informed these reprobates | I informed these reprobates that I WOULD NOT B... | |
| **3** | Greg Dunn | /users/62c35cdbacc0ea0012ccaffa | AU | 5 | 2024-09-17 07:15:49+00:00 | 1 | Advertise one price then increase it on website | I have bought from Amazon before and no proble... | |
| **4** | Sheila Hannah | /users/5ddbe429478d88251550610e | GB | 8 | 2024-09-16 18:37:17+00:00 | 1 | If I could give a lower rate I would | If I could give a lower rate I would! I cancel... | |
| **...** | ... | ... | ... | ... | ... | ... | ... | ... | |
| **21209** | Anders T | /users/47bd4ffe0000640001001044 | DK | 1 | 2009-03-22 13:14:12+00:00 | 5 | Fast!! | I have had perfect order fulfillment, and fast... | |
| **21210** | David E | /users/495bbbc0000064000100a972 | US | 2 | 2008-12-31 18:57:31+00:00 | 5 | Consistently Excellent | I have had perfect order fulfillment, and fast... | |
| **21211** | Joseph Harding | /users/48cfacbf0000640001005d04 | GB | 3 | 2008-09-16 13:05:05+00:00 | 3 | Good prices but delivery can take time : ( | I always find myself going back to amazon beco... | |
| **21212** | Mads Dørup | /users/474aaec70000640001000a44 | DK | 82 | 2008-04-28 11:09:05+00:00 | 5 | World-class online shopping | I have placed an abundance of orders with Amaz... | |
| **21213** | Kim Fuglsang Kramer | /users/46d1ed150000640001000051 | DK | 2 | 2007-08-27 17:25:01+00:00 | 4 | No title | those goods i've ordered by Amazon.com, have b... | |

21212 rows × 12 columns

Next steps: Generate code with `rev_df`    🔘 View recommended plots    New interactive sheet

```
1 rev_df['has_negative_title'] = rev_df['review_title'].str.contains("w ointed", case=False, na=False)
2 rev_df['has_positive_title'] = rev_df['review_title'].str.contains("great|excellent|amazing|perfect|love", case=False, na=False)
3 rev_df
```

| | reviewer_name | profile_link | country | review_count | review_date | rating | review_title | review_text | dat |
|---|---|---|---|---|---|---|---|---|---|
| 0 | Eugene ath | /users/66e8185ff1598352d6b3701a | US | 1 | 2024-09-16 13:44:26+00:00 | 1 | A Store That Doesn't Want to Sell Anything | I registered on the website, tried to order a ... | |
| 1 | Daniel ohalloran | /users/5d75e460200c1f6a6373648c | GB | 9 | 2024-09-16 18:26:46+00:00 | 1 | Had multiple orders one turned up and… | Had multiple orders one turned up and driver h... | |
| 2 | p fisher | /users/546cfcf1000064000197b88f | GB | 90 | 2024-09-16 21:47:39+00:00 | 1 | I informed these reprobates | I informed these reprobates that I WOULD NOT B... | |
| 3 | Greg Dunn | /users/62c35cdbacc0ea0012ccaffa | AU | 5 | 2024-09-17 07:15:49+00:00 | 1 | Advertise one price then increase it on website | I have bought from Amazon before and no proble... | |
| 4 | Sheila Hannah | /users/5ddbe429478d88251550610e | GB | 8 | 2024-09-16 18:37:17+00:00 | 1 | If I could give a lower rate I would | If I could give a lower rate I would! I cancel... | |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 21209 | Anders T | /users/47bd4ffe0000640001001044 | DK | 1 | 2009-03-22 13:14:12+00:00 | 5 | Fast!! | I have had perfect order fulfillment, and fast... | |
| 21210 | David E | /users/495bbbc0000064000100a972 | US | 2 | 2008-12-31 18:57:31+00:00 | 5 | Consistently Excellent | I have had perfect order fulfillment, and fast... | |
| 21211 | Joseph Harding | /users/48cfacbf0000640001005d04 | GB | 3 | 2008-09-16 13:05:05+00:00 | 3 | Good prices but delivery can take time : ( | I always find myself going back to amazon beco... | |
| 21212 | Mads Dørup | /users/474aaec70000640001000a44 | DK | 82 | 2008-04-28 11:09:05+00:00 | 5 | World-class online shopping | I have placed an abundance of orders with Amaz... | |
| 21213 | Kim Fuglsang Kramer | /users/46d1ed150000640001000051 | DK | 2 | 2007-08-27 17:25:01+00:00 | 4 | No title | those goods i've ordered by Amazon.com, have b... | |

21212 rows × 14 columns

Next steps: Generate code with `rev_df`    View recommended plots    New interactive sheet

```
1 # from google.colab import files
2 # rev_df.to_csv('cleaned_ecomm_data.csv', index=False)
3 # files.download('cleaned_ecomm_data.csv')
```
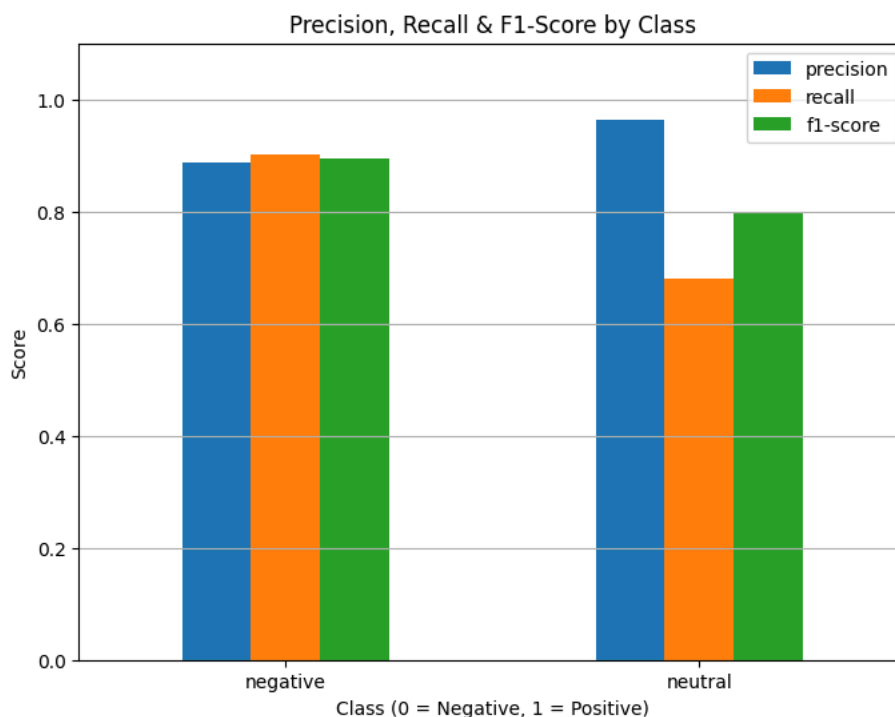
## Machine Learning

```
1 # Importing required libraries for prediction using machine learning
2
3 from sklearn.linear_model import LogisticRegression
4 from sklearn.feature_extraction.text import TfidfVectorizer
5 from sklearn.model_selection import train_test_split
6 from sklearn.metrics import classification_report, accuracy_score
7 import matplotlib.pyplot as plt
8
9 tf = TfidfVectorizer()
10 x = tf.fit_transform(rev_df['text_features'])
11 y = rev_df['sentiment']
12
13 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2, random_state=42)
14
```

```
15 model = LogisticRegression()
16 model.fit(x, y)
17 y_pred = model.predict(x_test)
18
19 print(classification_report(y_test, y_pred))
20 report = classification_report(y_test, y_pred, output_dict=True)
21 metrics_df = pd.DataFrame(report).transpose().iloc[:2][['precision', 'recall', 'f1-score']]
22 metrics_df.plot(kind='bar', figsize=(8, 6))
23 plt.title('Precision, Recall & F1-Score by Class')
24 plt.xlabel('Class (0 = Negative, 1 = Positive)')
25 plt.ylabel('Score')
26 plt.ylim(0, 1.1)
27 plt.grid(True, axis='y')
28 plt.xticks(rotation=0)
29 plt.show()
```

```
              precision    recall  f1-score   support

    negative       0.89      0.90      0.89      1692
     neutral       0.96      0.68      0.80       465
    positive       0.89      0.94      0.91      2086

    accuracy                           0.89      4243
   macro avg       0.91      0.84      0.87      4243
weighted avg       0.90      0.89      0.89      4243
```



Precision, Recall & F1-Score by Class

## ˅ Cross-Validation

```
 1 from sklearn.model_selection import cross_val_score, StratifiedKFold, cross_val_predict
 2 from sklearn.metrics import classification_report, accuracy_score
 3
 4 tf = TfidfVectorizer()
 5 x = tf.fit_transform(rev_df['text_features'])
 6 y = rev_df['sentiment']
 7
 8 model = LogisticRegression()
 9
10 skf = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
11 cv_scores = cross_val_score(model, x, y, cv=skf, scoring='accuracy')
12
13 print(f'Cross-Validation Scores : {cv_scores[0]}')
14
15 y_pred_cv = cross_val_predict(model, x, y, cv=skf)
16 print(f'Classification Report : {classification_report(y, y_pred_cv)}')
17 print(f'Accuracy Score : {accuracy_score(y, y_pred_cv):.2f}')
18
19 report = classification_report(y, y_pred_cv, output_dict=True)
20 metrics_df = pd.DataFrame(report).transpose().iloc[:2][['precision', 'recall', 'f1-score']]
21
22 # Plot metrics
```

```
23 metrics_df.plot(kind='bar', figsize=(8, 6))
24 plt.title('Precision, Recall & F1-Score by Class (Cross-Validated)')
25 plt.xlabel('Class (0 = Negative, 1 = Positive)')
26 plt.ylabel('Score')
27 plt.ylim(0, 1.1)
28 plt.grid(True, axis='y')
29 plt.xticks(rotation=0)
30 plt.tight_layout()
31 plt.show()
32
```
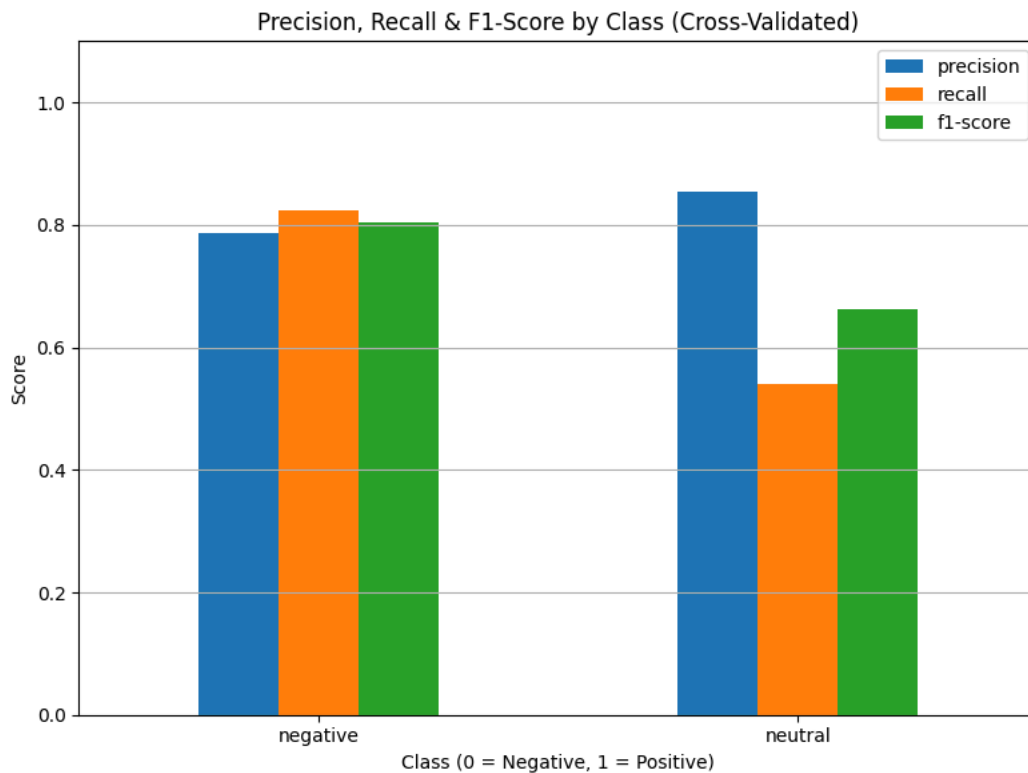
```
Cross-Validation Scores : 0.8197030403016733
Mean Accuracy : 0.81
Classification Report :              precision    recall  f1-score   support

       negative       0.79      0.82      0.80      8340
        neutral       0.85      0.54      0.66      2251
       positive       0.83      0.86      0.85     10621

       accuracy                           0.81     21212
      macro avg       0.82      0.74      0.77     21212
   weighted avg       0.81      0.81      0.81     21212

Accuracy Score : 0.81
```



Precision, Recall & F1-Score by Class (Cross-Validated)

## ˅ NLP

```
1 from textblob import TextBlob
2 import nltk
3 nltk.download('punkt_tab')
4
5 aspects = ['price', 'delivery', 'quality']
6 def aspect_sentiment(text):
7     blob = TextBlob(text)
8     aspect_sentiments = {}
9     for aspect in aspects:
10      for sentence in blob.sentences:
11        if aspect in sentence.lower():
12          aspect_sentiments[aspect] = sentence.sentiment.polarity
13    return aspect_sentiments
14
15 rev_df['Aspect_Sentiment'] = rev_df['text_features'].apply(aspect_sentiment)
16 print(rev_df[['text_features','Aspect_Sentiment']].head())
```

```
[nltk_data] Downloading package punkt_tab to /root/nltk_data...
[nltk_data]   Unzipping tokenizers/punkt_tab.zip.
                                       text_features    Aspect_Sentiment
0  I registered on the website, tried to order a ...               {}
1  Had multiple orders one turned up and driver h...  {'delivery': 0.0}
```