

```

In [1]: ▶ 1 import json
          2 import pandas as pd
          3
          4 data_file_path = 'C:/Users/ankit/yelp_merged_data.csv'
          5 checkin_df = pd.read_csv(data_file_path)
          6
          7 # Getting number of columns
          8 num_columns = checkin_df.shape[1]
          9
         10 print(f"The number of columns in the dataset is: {num_columns}")
         11 # Getting all the columns in the DataFrame
         12 print("Name of the Columns in the dataset:")
         13 for column in checkin_df.columns:
         14     print(column)
         15
         16 # Printing shape and details of the dataset
         17 print("Shape of the dataset:")
         18 print(checkin_df.shape)
         19
         20 # Printing Column Names:
         21 print("Column names:")
         22 print(checkin_df.columns)
         23
         24 # Printing Column Name with Data Type
         25 print("Datatype of each column:")
         26 print(checkin_df.dtypes)
         27
         28 # Printing first 5 rows of the dataset
         29 print("Few dataset entries:")
         30 print(checkin_df.head())
         31
         32 # Printing summart of the dataset
         33 checkin_df.describe(include='all')

```

```

0      Monday: "7:30-13:0", Tuesday: "7:30-13:0", ...
1  {"Monday": "6:30-20:30", "Tuesday": "6:30-20:3...
2                                     NaN
3  {"Tuesday": "11:0-21:0", "Wednesday": "11:0-21:0...
4  {"Monday": "0:0-0:0", "Friday": "11:0-17:0", "...

```

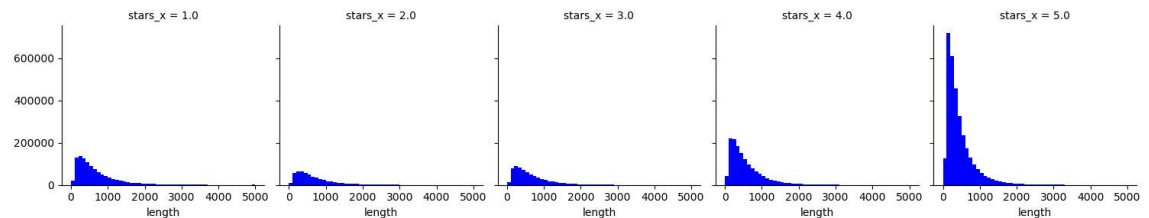
[5 rows x 21 columns]

Out[1]:

	review_id	user_id	business_id	st
count	6990280	6990280	6990280	6.990280
unique	6990280	1987929	150346	
top	KU_O5udG6zpxOg-VcAEodg	_BcWyKQL16ndpBdggh2kNA	_ab50qdWOk0DdB6XOrBitw	
freq	1	3048	7673	

```
In [3]: 1 import seaborn as sns
2 import matplotlib.pyplot as plt
3 #Adding and creating a new column named "Length" to count the number of
4 checkin_df['length'] = checkin_df['text'].apply(len)
5 checkin_df.head()
6
7 graph = sns.FacetGrid(data=checkin_df,col='stars_x')
8 graph.map(plt.hist,'length',bins=50,color='blue')
```

Out[3]: <seaborn.axisgrid.FacetGrid at 0x2154b05f410>

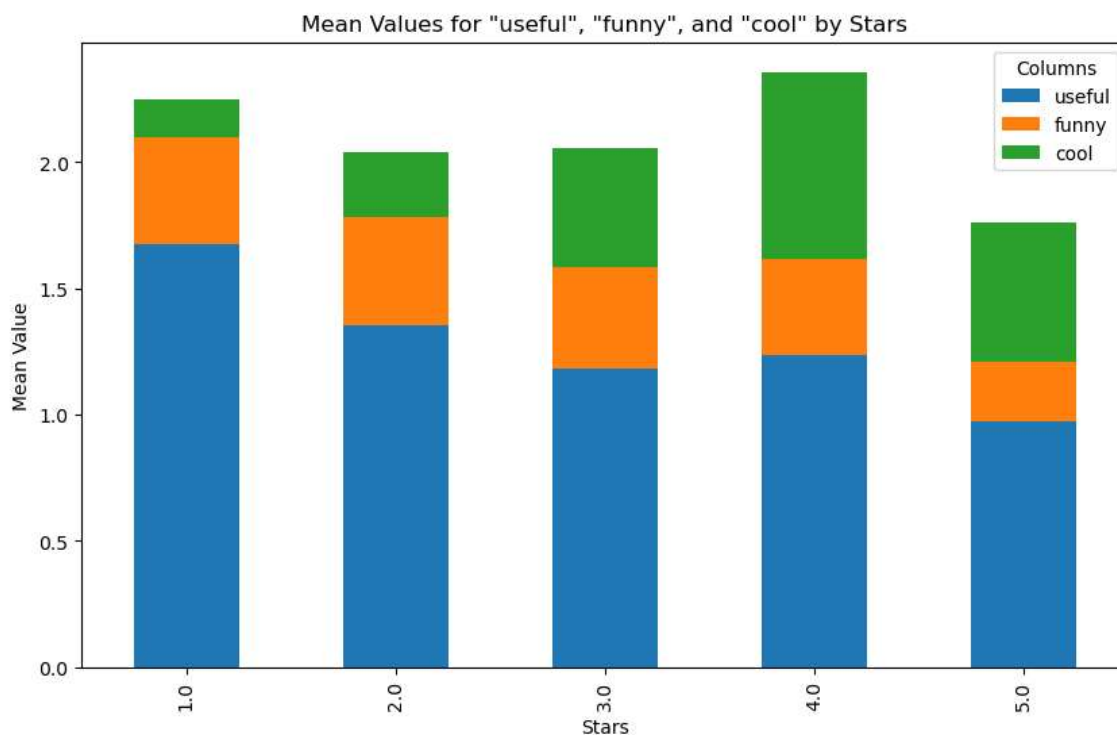


```
In [10]: 1 st_mean = checkin_df.groupby('stars_x').mean(numeric_only=True)
2 st_mean
```

Out[10]:

	useful	funny	cool	postal code	latitude	longitude	stars_y	review_count
stars_x								
1.0	1.673746	0.422610	0.154404	NaN	36.071437	-89.744251	3.009984	201.685561
2.0	1.354246	0.426679	0.258720	NaN	36.221857	-89.271493	3.434126	338.673870
3.0	1.181042	0.400245	0.472440	NaN	36.227045	-88.885623	3.610429	418.176550
4.0	1.234359	0.380260	0.742942	NaN	36.184931	-88.972196	3.799850	443.546784
5.0	0.972549	0.237988	0.548700	NaN	35.695603	-90.320310	4.058462	385.808362

```
In [11]: 1 import matplotlib.pyplot as plt
2
3 selected_columns = ['useful', 'funny', 'cool']
4 stval_selected = st_mean[selected_columns]
5
6 stval_selected.plot(kind='bar', stacked=True, figsize=(10, 6))
7 plt.title('Mean Values for "useful", "funny", and "cool" by Stars')
8 plt.xlabel('Stars')
9 plt.ylabel('Mean Value')
10 plt.legend(title='Columns', bbox_to_anchor=(1, 1))
11
12 plt.show()
```



In [12]:  1 st_mean.corr()

Out[12]:

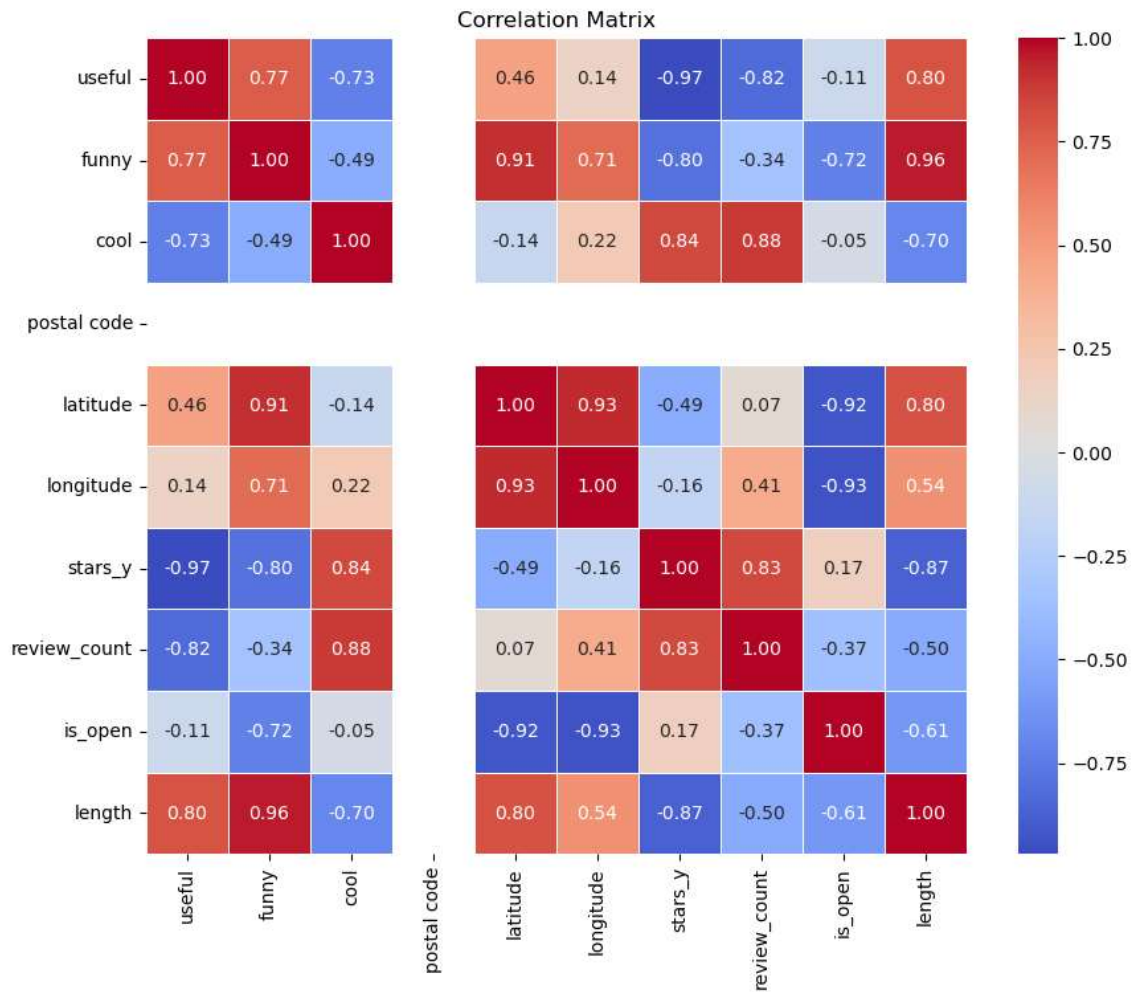
	useful	funny	cool	postal code	latitude	longitude	stars_y	review_
useful	1.000000	0.768755	-0.729256	NaN	0.458122	0.141436	-0.968761	-0.824346
funny	0.768755	1.000000	-0.485747	NaN	0.913415	0.705389	-0.796918	-0.715643
cool	-0.729256	-0.485747	1.000000	NaN	-0.135952	0.215737	0.835427	0.883743
postal code	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
latitude	0.458122	0.913415	-0.135952	NaN	1.000000	0.927746	-0.486565	0.066983
longitude	0.141436	0.705389	0.215737	NaN	0.927746	1.000000	-0.162000	0.408990
stars_y	-0.968761	-0.796918	0.835427	NaN	-0.486565	-0.162000	1.000000	0.833577
review_count	-0.824346	-0.344793	0.883743	NaN	0.066983	0.408990	0.833577	1.000000
is_open	-0.106978	-0.715643	-0.049289	NaN	-0.920341	-0.934607	0.171321	-0.049289
length	0.799799	0.960635	-0.695158	NaN	0.804144	0.539141	-0.874715	-0.695158



```

In [13]: ▶ 1 import seaborn as sns
2 import matplotlib.pyplot as plt
3
4 correlation_matrix = st_mean.corr()
5
6 # Creating a heatmap
7 plt.figure(figsize=(10, 8))
8 sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', fmt=".2f")
9 plt.title('Correlation Matrix')
10 plt.show()

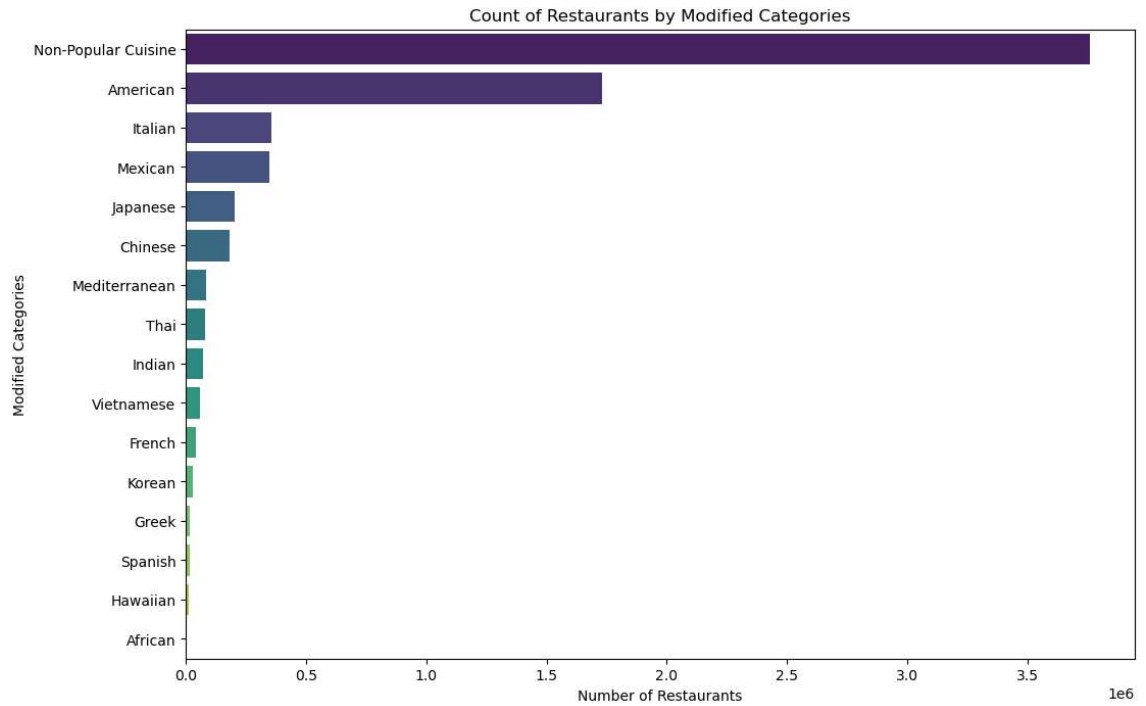
```



```

In [50]: ▶ 1 # Create the modified 'categories' column based on the mapping
2 checkin_df['modified_categories'] = checkin_df['categories'].apply(
3     lambda x: 'American' if 'restaurant' in x.lower() and 'american' in x
4         'Mexican' if 'restaurant' in x.lower() and 'mexican' in x
5         'Italian' if 'restaurant' in x.lower() and 'italian' in x
6         'Japanese' if 'restaurant' in x.lower() and 'japanese' in x
7         'Chinese' if 'restaurant' in x.lower() and 'chinese' in x
8         'Thai' if 'restaurant' in x.lower() and 'thai' in x.lower()
9         'Mediterranean' if 'restaurant' in x.lower() and 'mediterranean' in x
10        'French' if 'restaurant' in x.lower() and 'french' in x.lower()
11        'Vietnamese' if 'restaurant' in x.lower() and 'vietnamese' in x.lower()
12        'Greek' if 'restaurant' in x.lower() and 'greek' in x.lower()
13        'Indian' if 'restaurant' in x.lower() and 'indian' in x.lower()
14        'Korean' if 'restaurant' in x.lower() and 'korean' in x.lower()
15        'Hawaiian' if 'restaurant' in x.lower() and 'hawaiian' in x.lower()
16        'African' if 'restaurant' in x.lower() and 'african' in x.lower()
17        'Spanish' if 'restaurant' in x.lower() and 'spanish' in x.lower()
18        'Middle_eastern' if 'restaurant' in x.lower() and 'middle_eastern' in x.lower()
19        'Non-Popular Cuisine'
20 )
21
22
23 modified_categories_counts = checkin_df['modified_categories'].value_counts()
24
25 modified_categories_counts = modified_categories_counts.sort_values(ascending=False)
26
27 # Plot the graph of value counts
28 plt.figure(figsize=(12, 8))
29 sns.barplot(x=modified_categories_counts, y=modified_categories_counts)
30 plt.title('Count of Restaurants by Modified Categories')
31 plt.xlabel('Number of Restaurants')
32 plt.ylabel('Modified Categories')
33 plt.show()
34 # Count the categories in the 'modified_categories' column
35 #modified_categories_counts = checkin_df['modified_categories'].value_counts()
36
37 #for category, count in modified_categories_counts.items():
38     #print(f"{category}: {count} restaurants")
39

```



```
In [81]: 1 # CLASSIFICATION
2 checkin_df['review_category'] = pd.cut(checkin_df['stars_x'], bins=[0,
3
4 # Filter the DataFrame based on the new 'review_category' column
5 data_classes = checkin_df[checkin_df['review_category'].isin(['bad', 'a
6 data_classes.head()
7 print(data_classes.shape)
8
9 # Separate the dataset into X and Y for prediction
10 #x = data_classes['text'].iloc[:10000]
11 x = data_classes['text']
12 y = data_classes['review_category']
13 #y = data_classes['review_category'].iloc[:10000]
14
15 print(x.head())
16 print(y.head())
```

(6990280, 24)

```
0 If you decide to eat here, just be aware it is...
1 I've taken a lot of spin classes over the year...
2 Family diner. Had the buffet. Eclectic assortm...
3 Wow! Yummy, different, delicious. Our favo...
4 Cute interior and owner (?) gave us tour of up...
```

Name: text, dtype: object

```
0 average
1 good
2 average
3 good
4 good
```

Name: review_category, dtype: category

Categories (3, object): ['bad' < 'average' < 'good']

```
In [70]: 1 import nltk
2 from nltk.corpus import stopwords
3 import string
4 from nltk.tokenize import sent_tokenize
5 sent_tokenizer = sent_tokenize
6 from sklearn.feature_extraction.text import CountVectorizer
7
8 # CLEANING THE REVIEWS - REMOVAL OF STOPWORDS AND PUNCTUATION
9 def text_process(text):
10     nopunc = [char for char in text if char not in string.punctuation]
11     nopunc = ''.join(nopunc)
12     return [word for word in nopunc.split() if word.lower() not in stop
```

```
In [77]: 1 import nltk
2 from nltk.corpus import stopwords
3 import string
4 from nltk.tokenize import sent_tokenize
5 from sklearn.feature_extraction.text import CountVectorizer
6
7
8 x = checkin_df['text'].head(10000)
9 #x = checkin_df['text'] run this line of code for the whole data set, v
10 # but it is taking forever to complete preprocessing in our huge dataset
11
12 # CLEANING THE REVIEWS - REMOVAL OF STOPWORDS AND PUNCTUATION
13 def text_process(text):
14     nopunc = [char for char in text if char not in string.punctuation]
15     nopunc = ''.join(nopunc)
16     return [word for word in nopunc.split() if word.lower() not in stop
17
18 # CONVERTING THE WORDS INTO A VECTOR
19 vocab = CountVectorizer(analyzer=text_process).fit(x)
20 print(len(vocab.vocabulary_))
21
22 r0 = x.iloc[0]
23 vocab0 = vocab.transform([r0])
```

35071

```
In [78]: 1 # Transforming the features using the vocabulary
2 transformed_x = vocab.transform(x)
3
4 # Displaying information about the sparse matrix:
5 print("The shape of the sparse matrix is: ", transformed_x.shape)
6 print("Number of non-zero occurrences: ", transformed_x.nnz)
7
8 # density of the matrix:
9 density_percentage = (transformed_x.nnz / (transformed_x.shape[0] * tra
10 print("Density of the matrix: {:.2f}%".format(density_percentage))
```

The shape of the sparse matrix is: (10000, 35071)

Number of non-zero occurrences: 457423

Density of the matrix: 0.13%


```
In [92]: ▶ 1 from sklearn.model_selection import train_test_split
2 # SPLITTING THE DATASET INTO TRAINING SET AND TESTING SET
3 x_subset = x[:10000]
4 y_subset = y[:10000]
5 x_train, x_test, y_train, y_test = train_test_split(x_subset, y_subset,
```

```
In [90]: ▶ 1 vectorizer = CountVectorizer()
2
3 # Fit and transform the training data
4 x_train_vectorized = vectorizer.fit_transform(x_train)
5
6 # Transform the test data
7 x_test_vectorized = vectorizer.transform(x_test)
```

```

In [93]: ▶ 1 # Multinomial Naive Bayes
2 from sklearn.feature_extraction.text import TfidfVectorizer
3 from sklearn.model_selection import train_test_split
4 from sklearn.naive_bayes import MultinomialNB
5 from sklearn.metrics import confusion_matrix, accuracy_score, classification_report
6
7
8
9 tfidf_vectorizer = TfidfVectorizer()
10
11 # Fit and transform the training data
12 x_train_tfidf = tfidf_vectorizer.fit_transform(x_train)
13
14 # Transform the test data
15 x_test_tfidf = tfidf_vectorizer.transform(x_test)
16
17 # Creating and training the Multinomial Naive Bayes classifier
18 mnb_classifier = MultinomialNB()
19 mnb_classifier.fit(x_train_tfidf, y_train)
20
21 # Making predictions on the test set
22 predictions_mnb = mnb_classifier.predict(x_test_tfidf)
23
24 # Evaluating the Multinomial Naive Bayes classifier
25 print("Confusion Matrix for Multinomial Naive Bayes:")
26 print(confusion_matrix(y_test, predictions_mnb))
27
28 accuracy = accuracy_score(y_test, predictions_mnb) * 100
29 print("Accuracy Score: {:.2f}%".format(accuracy))
30
31 print("Classification Report:")
32 print(classification_report(y_test, predictions_mnb))
33

```

Confusion Matrix for Multinomial Naive Bayes:

```

[[ 0  0 213]
 [ 0  1 364]
 [ 0  0 1422]]

```

Accuracy Score: 71.15%

Classification Report:

	precision	recall	f1-score	support
average	0.00	0.00	0.00	213
bad	1.00	0.00	0.01	365
good	0.71	1.00	0.83	1422
accuracy			0.71	2000
macro avg	0.57	0.33	0.28	2000
weighted avg	0.69	0.71	0.59	2000

```
C:\Users\ankit\anaconda3\Lib\site-packages\sklearn\metrics\_classification.py:1469: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` parameter to control this behavior.
  _warn_prf(average, modifier, msg_start, len(result))
C:\Users\ankit\anaconda3\Lib\site-packages\sklearn\metrics\_classification.py:1469: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` parameter to control this behavior.
  _warn_prf(average, modifier, msg_start, len(result))
C:\Users\ankit\anaconda3\Lib\site-packages\sklearn\metrics\_classification.py:1469: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` parameter to control this behavior.
  _warn_prf(average, modifier, msg_start, len(result))
```

```
In [95]: ▶ 1 #Random Forest
2 from sklearn.ensemble import RandomForestClassifier
3 rf_classifier = RandomForestClassifier(n_estimators=100, random_state=1)
4 rf_classifier.fit(x_train_tfidf, y_train)
5
6
7 predictions_rf = rf_classifier.predict(x_test_tfidf)
8
9 # Evaluating the Random Forest Classifier
10 print("Confusion Matrix for Random Forest Classifier:")
11 print(confusion_matrix(y_test, predictions_rf))
12
13 accuracy = accuracy_score(y_test, predictions_rf) * 100
14 print("Accuracy Score: {:.2f}%".format(accuracy))
15
16 print("Classification Report:")
17 print(classification_report(y_test, predictions_rf))
```

Confusion Matrix for Random Forest Classifier:

```
[[ 1  10 202]
 [ 0 113 252]
 [ 0   1 1421]]
```

Accuracy Score: 76.75%

Classification Report:

	precision	recall	f1-score	support
average	1.00	0.00	0.01	213
bad	0.91	0.31	0.46	365
good	0.76	1.00	0.86	1422
accuracy			0.77	2000
macro avg	0.89	0.44	0.44	2000
weighted avg	0.81	0.77	0.70	2000

```
In [96]: ▶ 1 # Decision Tree
2 from sklearn.tree import DecisionTreeClassifier
3 dt_classifier = DecisionTreeClassifier(random_state=101)
4 dt_classifier.fit(x_train_tfidf, y_train)
5
6
7 predictions_dt = dt_classifier.predict(x_test_tfidf)
8
9
10 print("Confusion Matrix for Decision Tree Classifier:")
11 print(confusion_matrix(y_test, predictions_dt))
12
13 accuracy = accuracy_score(y_test, predictions_dt) * 100
14 print("Accuracy Score: {:.2f}%".format(accuracy))
15
16 print("Classification Report:")
17 print(classification_report(y_test, predictions_dt))
```

Confusion Matrix for Decision Tree Classifier:

```
[[ 53  38 122]
 [ 56 184 125]
 [136 148 1138]]
```

Accuracy Score: 68.75%

Classification Report:

	precision	recall	f1-score	support
average	0.22	0.25	0.23	213
bad	0.50	0.50	0.50	365
good	0.82	0.80	0.81	1422
accuracy			0.69	2000
macro avg	0.51	0.52	0.51	2000
weighted avg	0.70	0.69	0.69	2000

```
In [97]: ▶ 1 # Support Vector Machine
2 from sklearn.svm import SVC
3 svm_classifier = SVC(kernel='linear', random_state=101)
4 svm_classifier.fit(x_train_tfidf, y_train)
5
6
7 predictions_svm = svm_classifier.predict(x_test_tfidf)
8
9 # Evaluating the Support Vector Machine (SVM) Classifier
10 print("Confusion Matrix for Support Vector Machine (SVM) Classifier:")
11 print(confusion_matrix(y_test, predictions_svm))
12
13 accuracy = accuracy_score(y_test, predictions_svm) * 100
14 print("Accuracy Score: {:.2f}%".format(accuracy))
15
16 print("Classification Report:")
17 print(classification_report(y_test, predictions_svm))
```

Confusion Matrix for Support Vector Machine (SVM) Classifier:

```
[[ 41  43 129]
 [ 18 289  58]
 [ 27  16 1379]]
```

Accuracy Score: 85.45%

Classification Report:

	precision	recall	f1-score	support
average	0.48	0.19	0.27	213
bad	0.83	0.79	0.81	365
good	0.88	0.97	0.92	1422
accuracy			0.85	2000
macro avg	0.73	0.65	0.67	2000
weighted avg	0.83	0.85	0.83	2000

```
In [100]: ▶ 1 # K Nearest Neighbour Algorithm
2 from sklearn.neighbors import KNeighborsClassifier
3 knn_classifier = KNeighborsClassifier(n_neighbors=5)
4 knn_classifier.fit(x_train_tfidf, y_train)
5
6 # Making predictions on the test set
7 predictions_knn = knn_classifier.predict(x_test_tfidf)
8
9
10 print("Confusion Matrix for k-Nearest Neighbors (KNN) Classifier:")
11 print(confusion_matrix(y_test, predictions_knn))
12
13 accuracy = accuracy_score(y_test, predictions_knn) * 100
14 print("Accuracy Score: {:.2f}%".format(accuracy))
15
16 print("Classification Report:")
17 print(classification_report(y_test, predictions_knn))
```

Confusion Matrix for k-Nearest Neighbors (KNN) Classifier:

```
[[ 40  51 122]
 [ 60 200 105]
 [119 140 1163]]
```

Accuracy Score: 70.15%

Classification Report:

	precision	recall	f1-score	support
average	0.18	0.19	0.19	213
bad	0.51	0.55	0.53	365
good	0.84	0.82	0.83	1422
accuracy			0.70	2000
macro avg	0.51	0.52	0.51	2000
weighted avg	0.71	0.70	0.70	2000

```

In [105]: ▶ 1 # XGBoost Classifier
2 #!pip install xgboost
3 import xgboost as xgb
4 from sklearn.model_selection import train_test_split
5 from sklearn.metrics import mean_squared_error
6 from xgboost import XGBClassifier
7 from sklearn.preprocessing import LabelEncoder
8
9 label_encoder = LabelEncoder()
10 y_train_encoded = label_encoder.fit_transform(y_train)
11 y_test_encoded = label_encoder.transform(y_test)
12
13
14 # Creating and training the XGBoost Classifier
15 xgb_classifier = xgb.XGBClassifier(
16     learning_rate=0.1,
17     max_depth=5,
18     subsample=0.8,
19     colsample_bytree=0.8,
20     n_estimators=100,
21     random_state=101
22 )
23 xgb_classifier.fit(x_train_tfidf, y_train_encoded)
24
25
26 predictions_xgb = xgb_classifier.predict(x_test_tfidf)
27 predictions_xgb_labels = label_encoder.inverse_transform(predictions_xgb)
28
29 # Evaluating the XGBoost Classifier
30 print("Confusion Matrix for XGBoost Classifier:")
31 print(confusion_matrix(y_test, predictions_xgb_labels))
32
33 accuracy = accuracy_score(y_test, predictions_xgb_labels) * 100
34 print("Accuracy Score: {:.2f}%".format(accuracy))
35
36 print("Classification Report:")
37 print(classification_report(y_test, predictions_xgb_labels))

```

Confusion Matrix for XGBoost Classifier:

```

[[ 27  28 158]
 [ 18 226 121]
 [ 16  20 1386]]

```

Accuracy Score: 81.95%

Classification Report:

	precision	recall	f1-score	support
average	0.44	0.13	0.20	213
bad	0.82	0.62	0.71	365
good	0.83	0.97	0.90	1422
accuracy			0.82	2000
macro avg	0.70	0.57	0.60	2000
weighted avg	0.79	0.82	0.79	2000

```

In [107]: ► 1 # MULTILAYER PERCEPTRON CLASSIFIER
2 from sklearn.neural_network import MLPClassifier
3
4 mlp_classifier = MLPClassifier(
5     hidden_layer_sizes=(100,),
6     activation='relu',
7     solver='adam',
8     max_iter=200,
9     random_state=101
10 )
11 mlp_classifier.fit(x_train_tfidf, y_train_encoded)
12
13 # Making predictions on the test set
14 predictions_mlp = mlp_classifier.predict(x_test_tfidf)
15
16 # Inverse transform the numerical labels back to original string labels
17 predictions_mlp_labels = label_encoder.inverse_transform(predictions_mlp)
18
19 # Evaluating the MLP Classifier
20 print("Confusion Matrix for MLP Classifier:")
21 print(confusion_matrix(y_test, predictions_mlp_labels))
22
23 accuracy = accuracy_score(y_test, predictions_mlp_labels) * 100
24 print("Accuracy Score: {:.2f}%".format(accuracy))
25
26 print("Classification Report:")
27 print(classification_report(y_test, predictions_mlp_labels))

```

Confusion Matrix for MLP Classifier:

```

[[ 58  47 108]
 [ 30 279  56]
 [ 58  32 1332]]

```

Accuracy Score: 83.45%

Classification Report:

	precision	recall	f1-score	support
average	0.40	0.27	0.32	213
bad	0.78	0.76	0.77	365
good	0.89	0.94	0.91	1422
accuracy			0.83	2000
macro avg	0.69	0.66	0.67	2000
weighted avg	0.82	0.83	0.82	2000


```

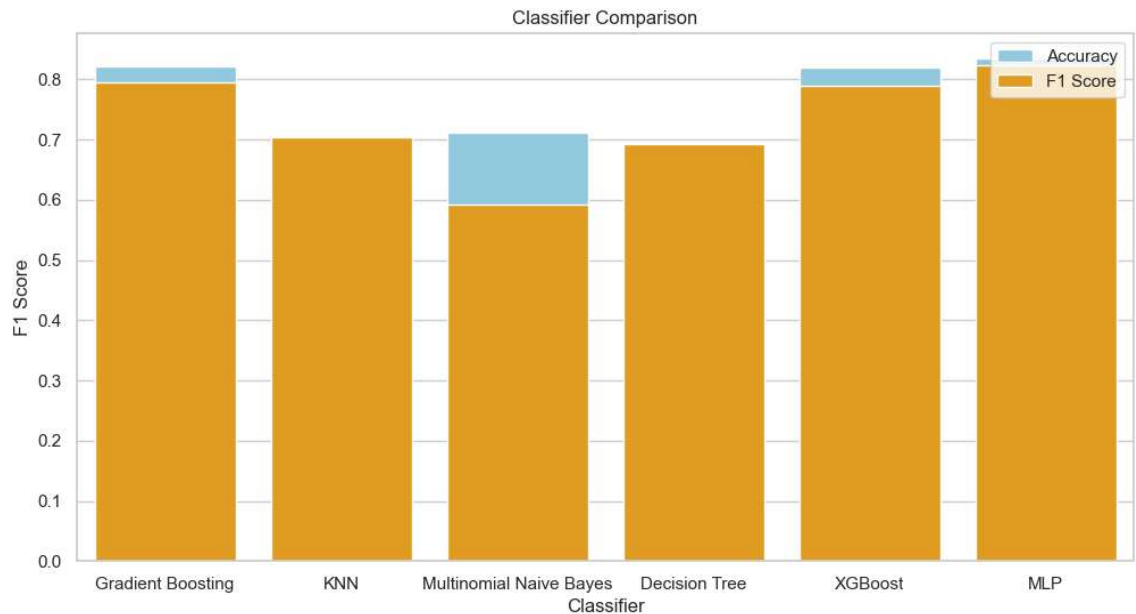
In [112]: 1 from sklearn.metrics import f1_score
2 from sklearn.metrics import accuracy_score, f1_score
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 import pandas as pd
6
7 #graph for results comparison
8 results = {'Classifier': [], 'Accuracy': [], 'F1 Score': []}
9
10 # Gradient Boosting Classifier
11 gbi = GradientBoostingClassifier(learning_rate=0.1, max_depth=5, max_fe
12 gbi.fit(x_train_tfidf, y_train)
13 pred_gbi = gbi.predict(x_test_tfidf)
14 accuracy_gbi = accuracy_score(y_test, pred_gbi)
15 f1_gbi = f1_score(y_test, pred_gbi, average='weighted')
16 results['Classifier'].append('Gradient Boosting')
17 results['Accuracy'].append(accuracy_gbi)
18 results['F1 Score'].append(f1_gbi)
19
20 # K-Nearest Neighbors (KNN) Classifier
21 knn_classifier = KNeighborsClassifier(n_neighbors=5)
22 knn_classifier.fit(x_train_tfidf, y_train_encoded)
23 pred_knn = knn_classifier.predict(x_test_tfidf)
24 accuracy_knn = accuracy_score(y_test_encoded, pred_knn)
25 f1_knn = f1_score(y_test_encoded, pred_knn, average='weighted')
26 results['Classifier'].append('KNN')
27 results['Accuracy'].append(accuracy_knn)
28 results['F1 Score'].append(f1_knn)
29
30 # Multinomial Naive Bayes Classifier
31 mnb_classifier = MultinomialNB()
32 mnb_classifier.fit(x_train_tfidf, y_train_encoded)
33 pred_mnb = mnb_classifier.predict(x_test_tfidf)
34 accuracy_mnb = accuracy_score(y_test_encoded, pred_mnb)
35 f1_mnb = f1_score(y_test_encoded, pred_mnb, average='weighted')
36 results['Classifier'].append('Multinomial Naive Bayes')
37 results['Accuracy'].append(accuracy_mnb)
38 results['F1 Score'].append(f1_mnb)
39
40 # Decision Tree Classifier
41 dt_classifier = DecisionTreeClassifier(random_state=101)
42 dt_classifier.fit(x_train_tfidf, y_train_encoded)
43 pred_dt = dt_classifier.predict(x_test_tfidf)
44 accuracy_dt = accuracy_score(y_test_encoded, pred_dt)
45 f1_dt = f1_score(y_test_encoded, pred_dt, average='weighted')
46 results['Classifier'].append('Decision Tree')
47 results['Accuracy'].append(accuracy_dt)
48 results['F1 Score'].append(f1_dt)
49
50 # XGBoost Classifier
51 xgb_classifier = xgb.XGBClassifier(learning_rate=0.1, max_depth=5, subs
52 xgb_classifier.fit(x_train_tfidf, y_train_encoded)
53 pred_xgb = xgb_classifier.predict(x_test_tfidf)
54 accuracy_xgb = accuracy_score(y_test_encoded, pred_xgb)
55 f1_xgb = f1_score(y_test_encoded, pred_xgb, average='weighted')
56 results['Classifier'].append('XGBoost')
57 results['Accuracy'].append(accuracy_xgb)

```

```

58 results['F1 Score'].append(f1_xgb)
59
60 # Multi-Layer Perceptron (MLP) Classifier
61 mlp_classifier = MLPClassifier(hidden_layer_sizes=(100,), activation='tanh')
62 mlp_classifier.fit(x_train_tfidf, y_train_encoded)
63 pred_mlp = mlp_classifier.predict(x_test_tfidf)
64 accuracy_mlp = accuracy_score(y_test_encoded, pred_mlp)
65 f1_mlp = f1_score(y_test_encoded, pred_mlp, average='weighted')
66 results['Classifier'].append('MLP')
67 results['Accuracy'].append(accuracy_mlp)
68 results['F1 Score'].append(f1_mlp)
69
70 # Convert results to DataFrame for easy plotting
71 import pandas as pd
72 df_results = pd.DataFrame(results)
73
74 # Plotting the results
75 plt.figure(figsize=(12, 6))
76 sns.set(style="whitegrid")
77 sns.barplot(x='Classifier', y='Accuracy', data=df_results, color='skyblue')
78 sns.barplot(x='Classifier', y='F1 Score', data=df_results, color='orange')
79 plt.title('Classifier Comparison')
80 plt.legend(loc='upper right')
81 plt.show()
82

```



In []: ▶

1