

Exploration

```
In [ ]: import pandas as pd
import matplotlib.pyplot as plt

trainingSet = pd.read_csv("../data/train.csv")
testingSet = pd.read_csv("../data/test.csv")

print("train.csv shape is ", trainingSet.shape)
print("test.csv shape is ", testingSet.shape)

print()

print(trainingSet.head())
print()
print(testingSet.head())

print()

print(trainingSet.describe())

trainingSet['Score'].value_counts().plot(kind='bar', legend=True, alpha=.5)
plt.title("Count of Scores")
plt.show()

trainingSet['ProductId'].value_counts().nlargest(25).plot(kind='bar', legend=True,
plt.title("Top 25 most rated Products")
plt.show()

trainingSet['ProductId'].value_counts().nsmallest(25).plot(kind='bar', legend=True,
plt.title("Top 25 least rated Products")
plt.show()

trainingSet['UserId'].value_counts().nlargest(25).plot(kind='bar', legend=True, alp
plt.title("Top 25 Reviewers")
plt.show()

trainingSet['UserId'].value_counts().nsmallest(25).plot(kind='bar', legend=True, al
plt.title("Lowest 25 Reviewers")
plt.show()

trainingSet[['Score', 'HelpfulnessNumerator']].groupby('Score').mean().plot(kind='ba
plt.title("Mean Helpfulness Numerator per Score")
plt.show()

trainingSet[['Score', 'ProductId']].groupby('ProductId').mean().nlargest(25, 'Score'
plt.title("Top 25 best rated Products")
plt.show()

trainingSet[['Score', 'ProductId']].groupby('ProductId').mean().nsmallest(25, 'Score
plt.title("Top 25 worst rated Products")
plt.show()

trainingSet[['Score', 'UserId']].groupby('UserId').mean().nlargest(25, 'Score').plot
plt.title("Top 25 kindest Reviewers")
plt.show()

trainingSet[['Score', 'UserId']].groupby('UserId').mean().nsmallest(25, 'Score').plo
plt.title("Top 25 harshest Reviewers")
plt.show()
```

```
trainingSet[trainingSet['ProductId'].isin(trainingSet['ProductId'].value_counts().nl
plt.title("Mean of top 25 most rated Products")
plt.show()
```

```
train.csv shape is (139753, 9)
test.csv shape is (17470, 2)
```

	Id	ProductId	UserId	HelpfulnessNumerator	\
0	195370	1890228583	A3VLX5Z090RQOV	1	
1	1632470	B00BEIYSL4	AUDXDMFM49NGY	0	
2	9771	0767809335	A3LFIA97BUU5IE	3	
3	218855	6300215792	A1QZM75342ZQVQ	1	
4	936225	B000B5X0ZW	ANM2SCEUL3WL1	1	

	HelpfulnessDenominator	Time	\
0	2	1030838400	
1	1	1405036800	
2	36	983750400	
3	1	1394841600	
4	1	1163721600	

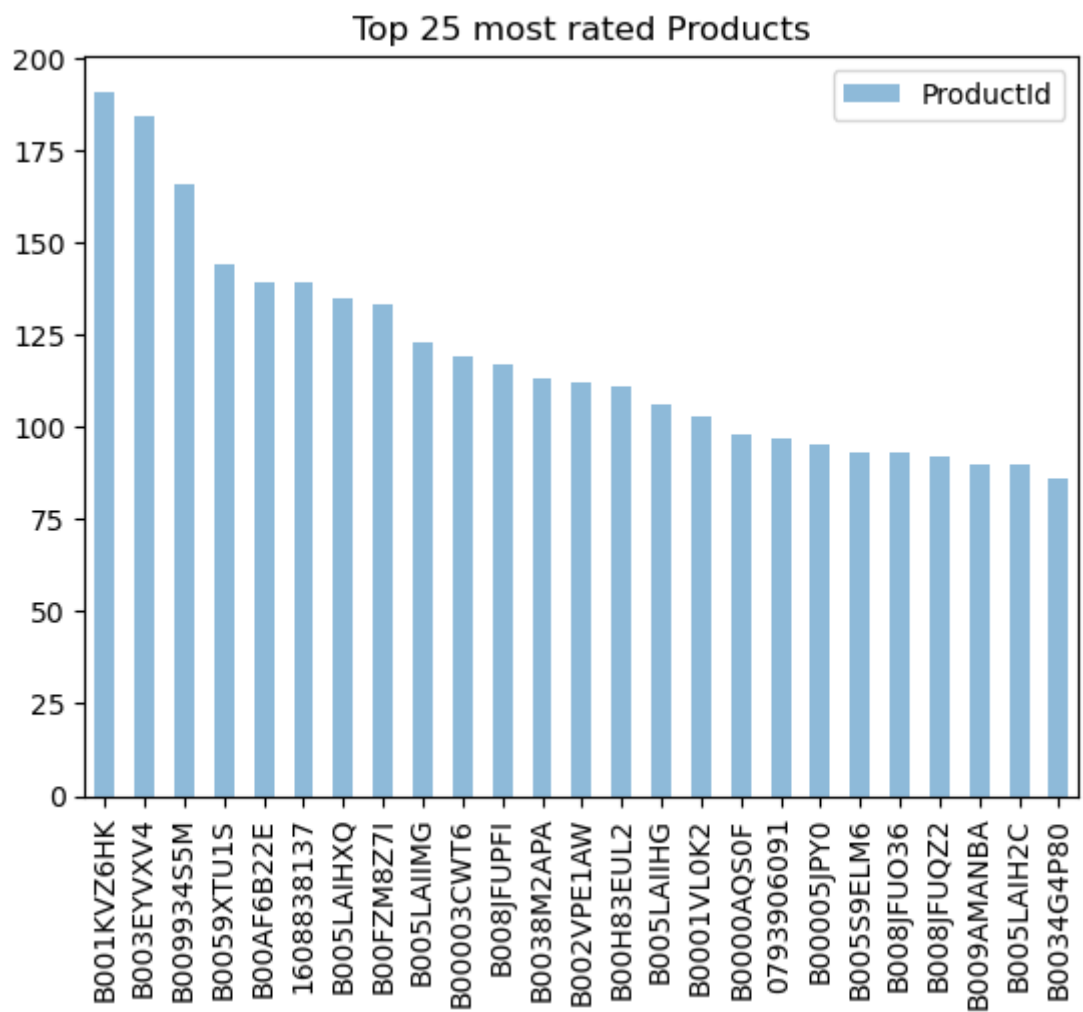
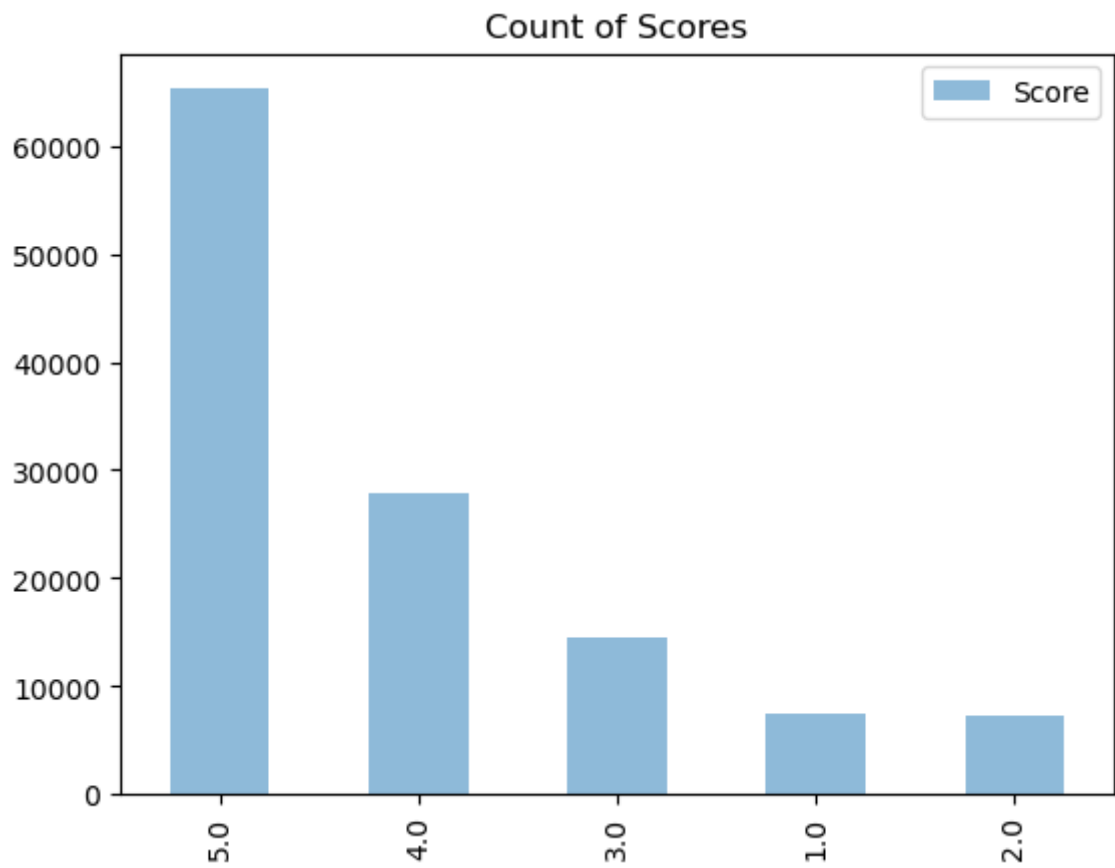
	Summary	\
0	An Unexplained Anime Review	
1	not great.	
2	Technical problem with this DVD	
3	Heeeeyyyyy LAAAADEEEE!!!!	
4	Herzog the Great Traveler of both natural and ...	

	Text	Score
0	I was very anxious to see the Uncut version of...	2.0
1	Movie was okay...not great.	3.0
2	Like the Dinosaur Collector's Edition DVD, thi...	1.0
3	Come on, now..... this has to be, by far, the...	5.0
4	I've always been a great admirer of Herzog's o...	4.0

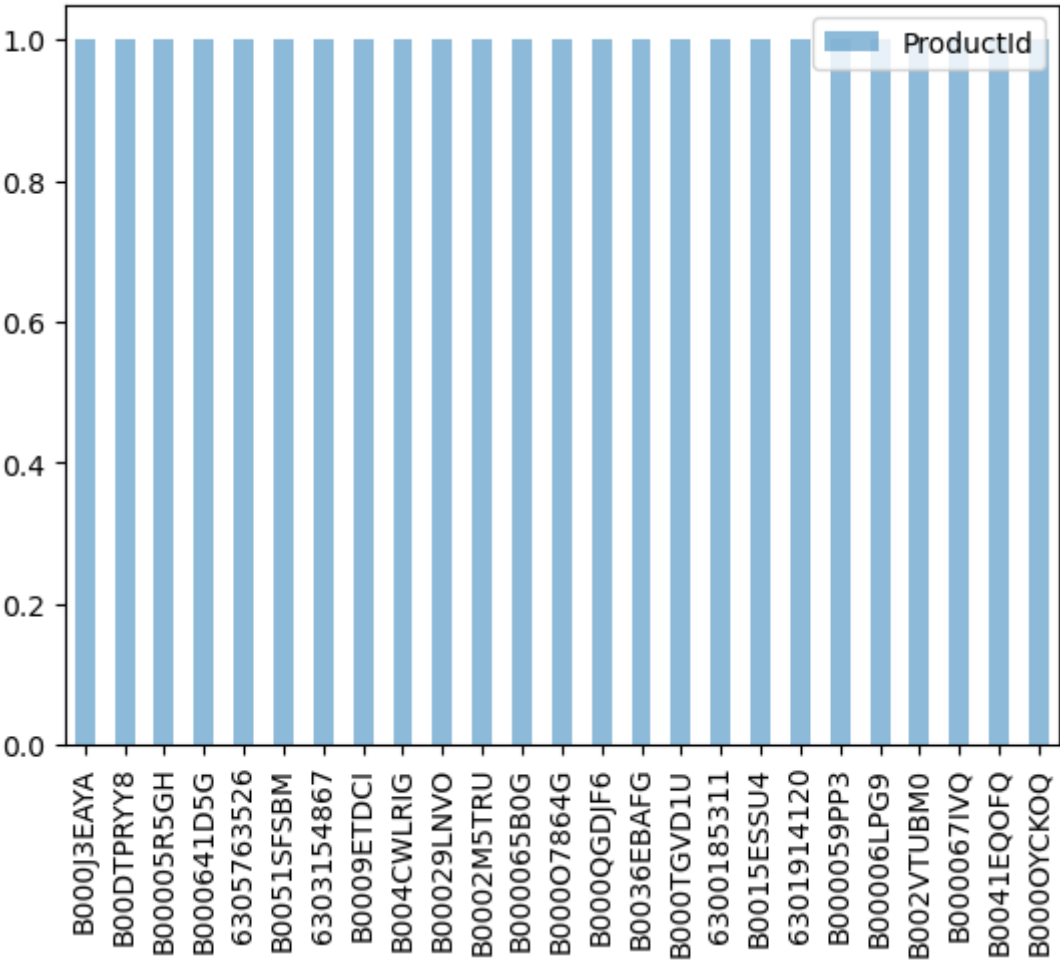
	Id	Score
0	786781	NaN
1	17153	NaN
2	1557328	NaN
3	1242666	NaN
4	1359242	NaN

	Id	HelpfulnessNumerator	HelpfulnessDenominator	\
count	1.397530e+05	139753.000000	139753.000000	
mean	8.497881e+05	3.601096	5.313246	
std	4.896942e+05	20.101195	22.300962	
min	8.000000e+00	0.000000	0.000000	
25%	4.258660e+05	0.000000	0.000000	
50%	8.510200e+05	1.000000	1.000000	
75%	1.273392e+06	3.000000	5.000000	
max	1.697519e+06	4646.000000	4682.000000	

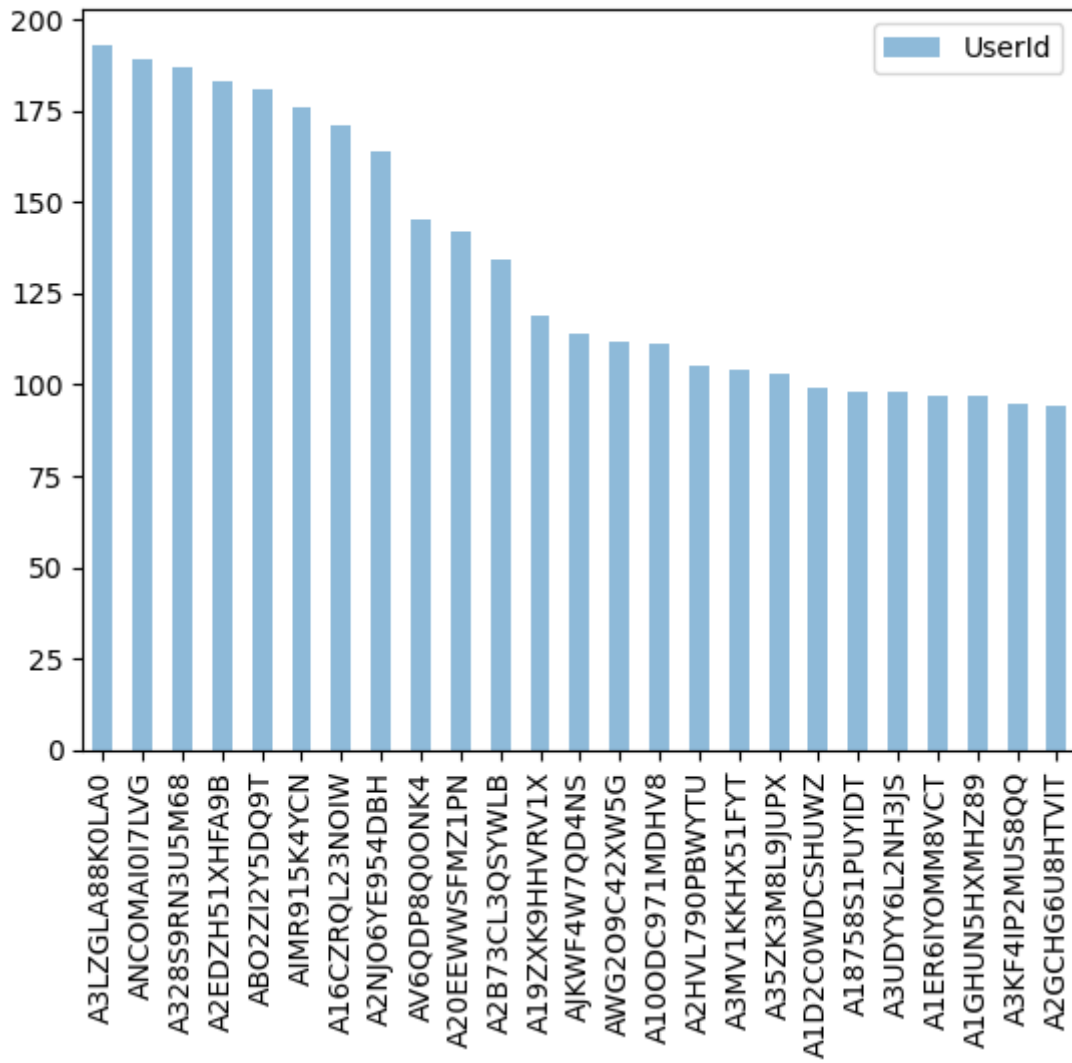
	Time	Score
count	1.397530e+05	122283.000000
mean	1.262516e+09	4.115552
std	1.287262e+08	1.191661
min	8.948448e+08	1.000000
25%	1.164758e+09	4.000000
50%	1.307318e+09	5.000000
75%	1.373155e+09	5.000000
max	1.406074e+09	5.000000



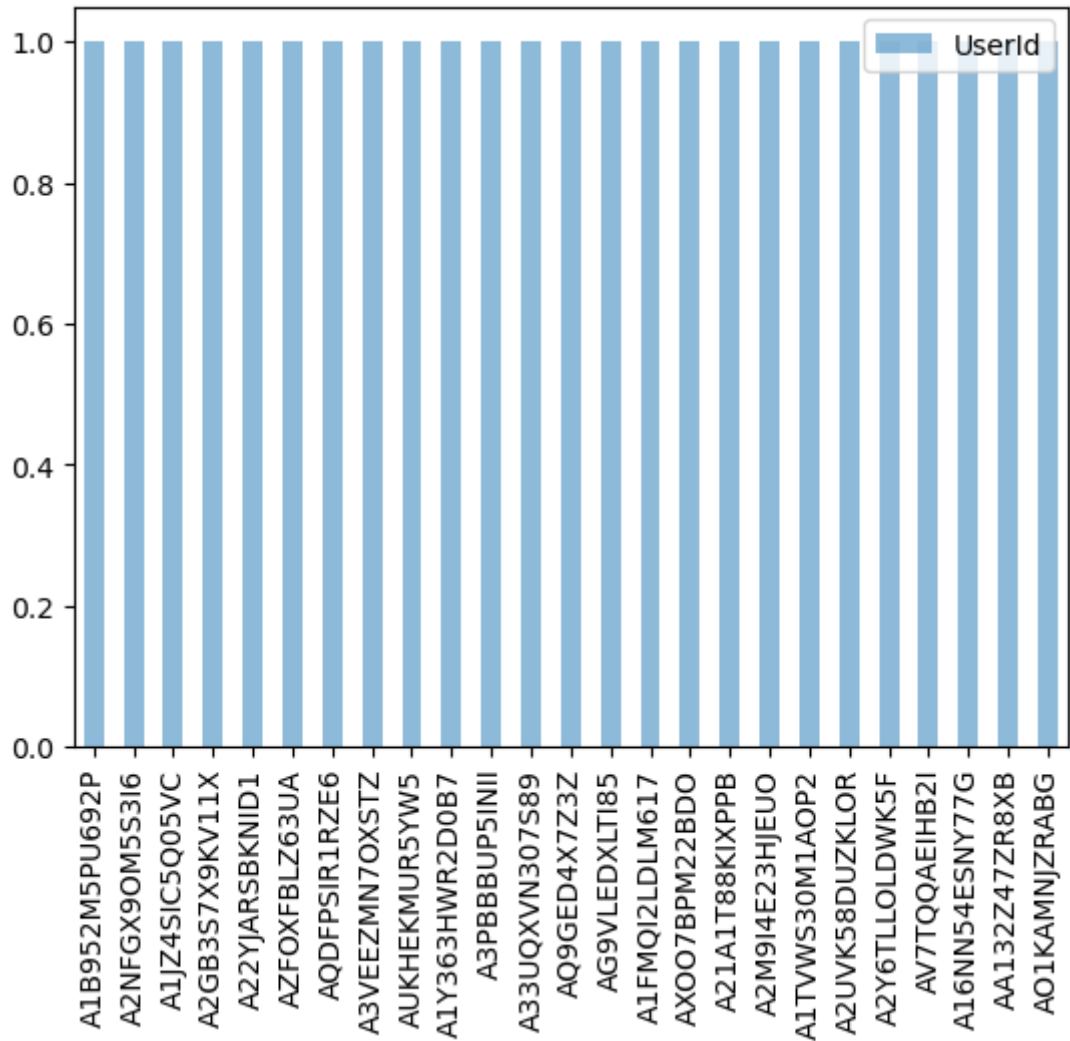
Top 25 least rated Products



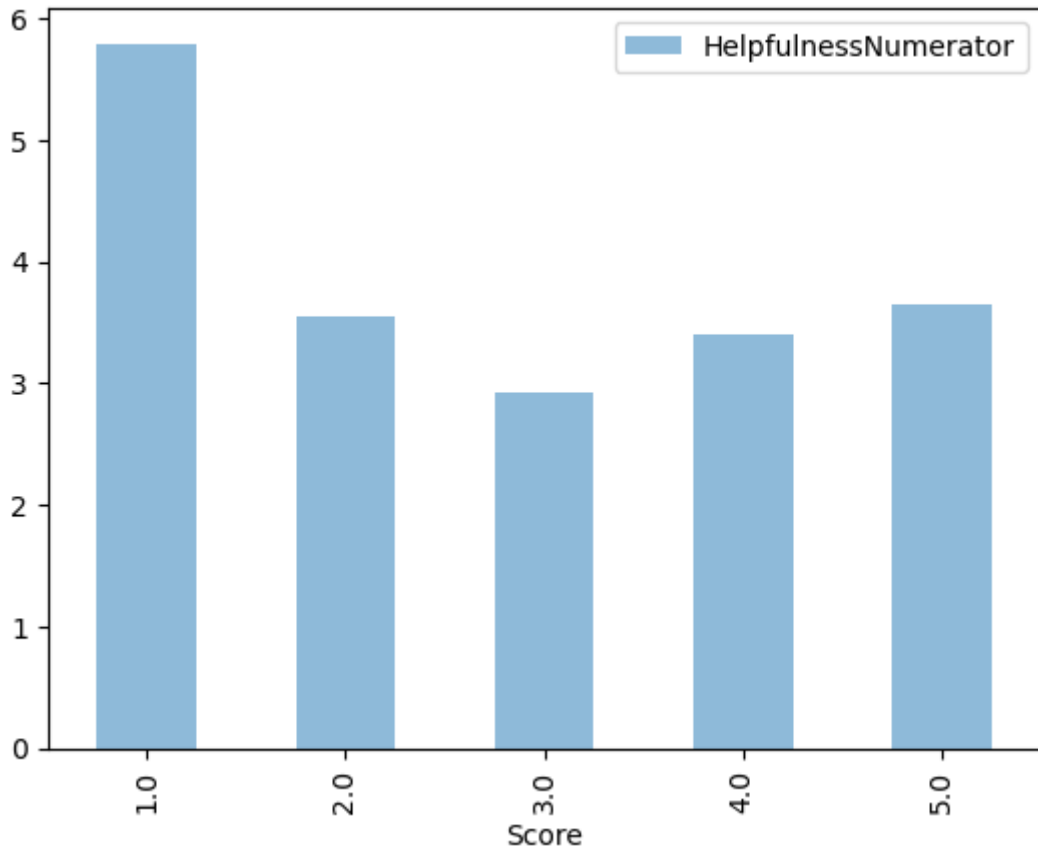
Top 25 Reviewers



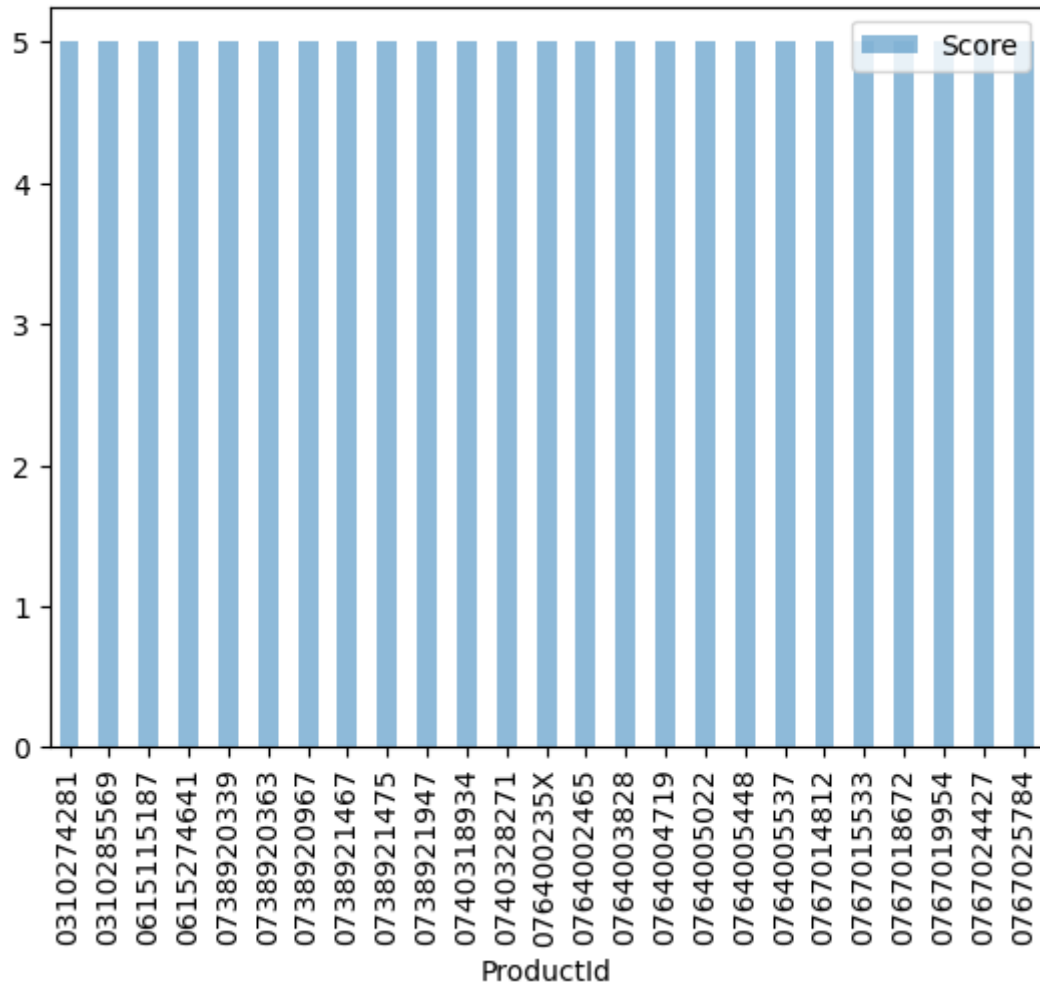
Lowest 25 Reviewers



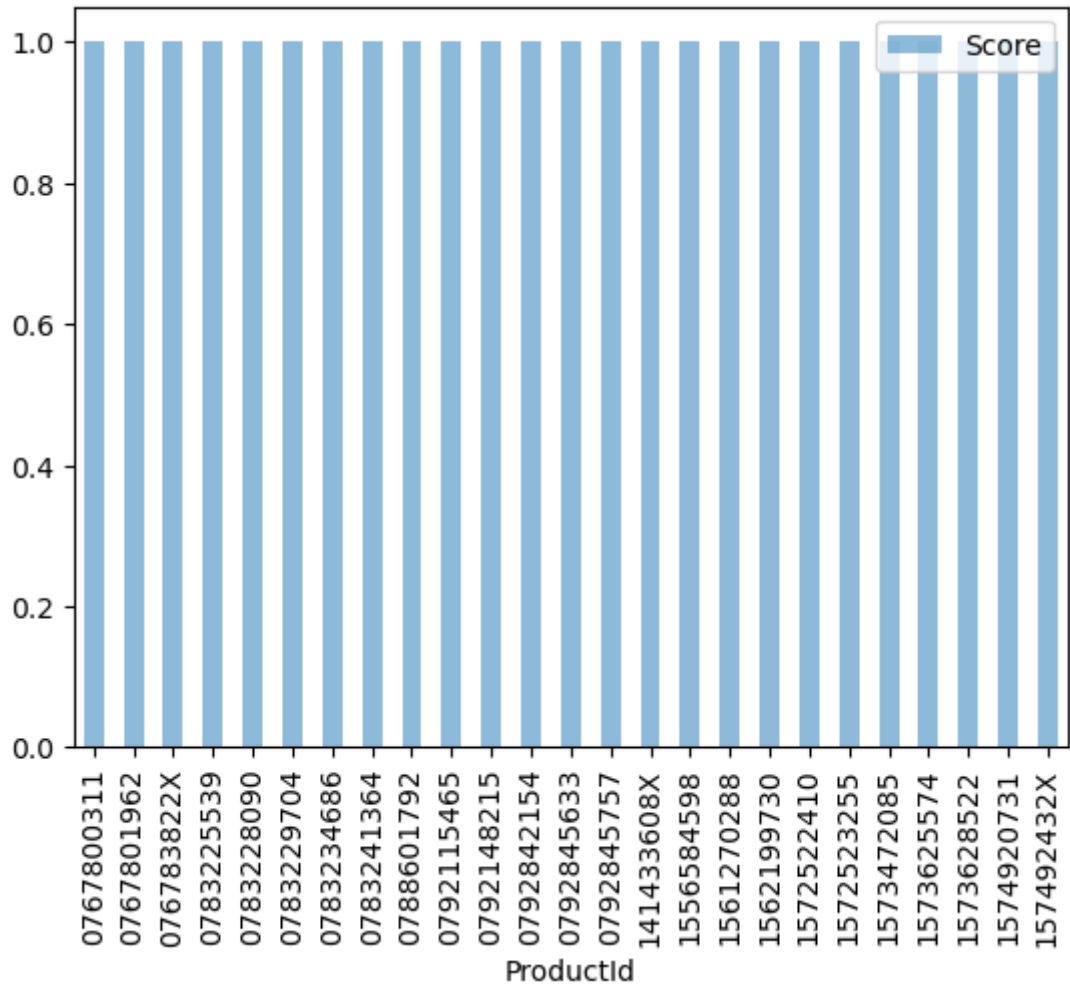
Mean Helpfulness Numerator per Score



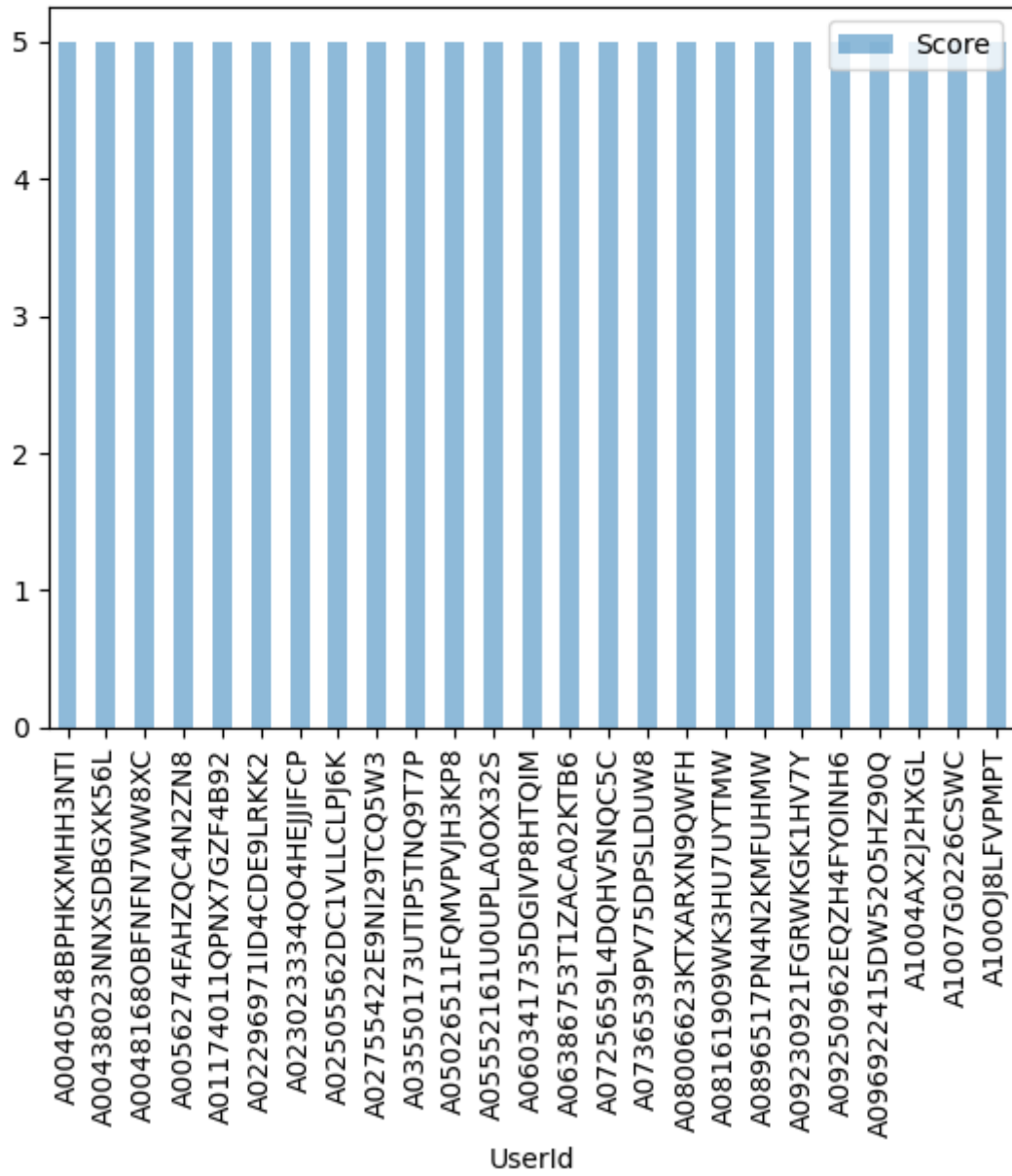
Top 25 best rated Products



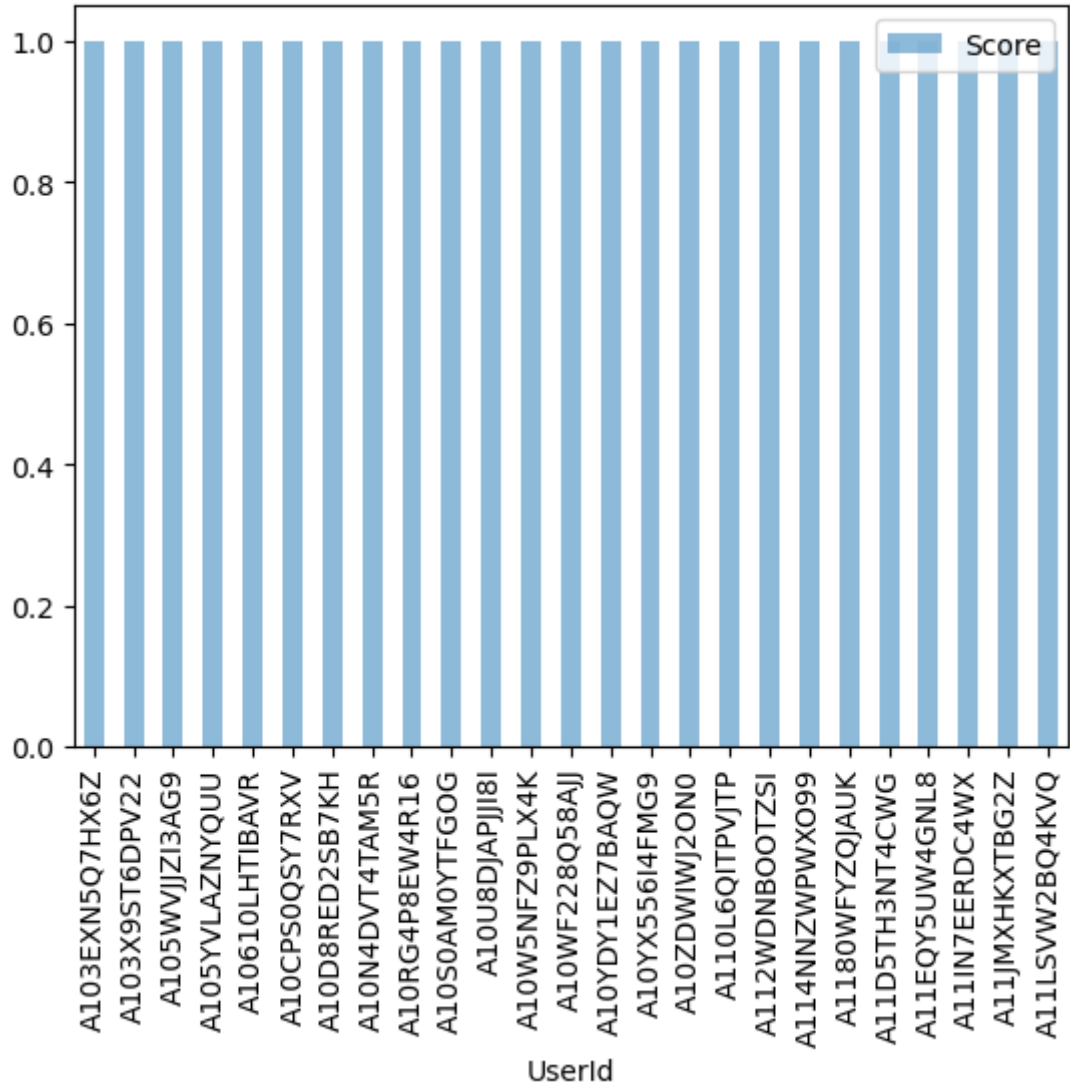
Top 25 worst rated Products

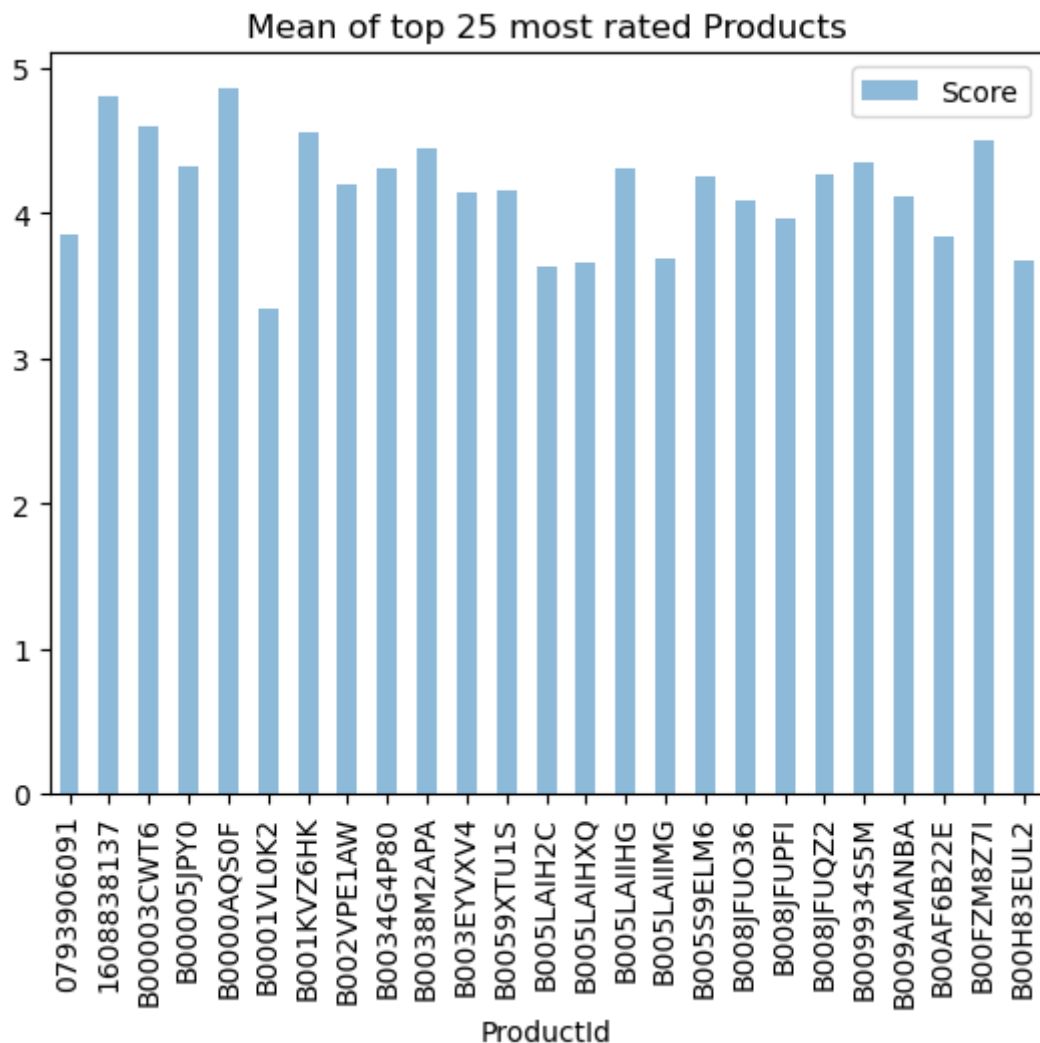


Top 25 kindest Reviewers



Top 25 harshest Reviewers





ctrl + / == #mutipal rows

take less score 5(5 is far more than other)

Feature Extraction

```
In [ ]: import pandas as pd
from sklearn.preprocessing import LabelEncoder
from sklearn.preprocessing import LabelEncoder, OneHotEncoder
from textblob import TextBlob

def process(df):
    # This is where you can do all your processing

    #this part is just for test,the code used for train is outside this function

    #mean_denominator = df[df['HelpfulnessDenominator'] != 0]['HelpfulnessDenominator']
    #df.loc[df['HelpfulnessDenominator'] == 0, 'HelpfulnessDenominator'] = mean_denominator
    df['Helpfulness'] = df['HelpfulnessNumerator'] / df['HelpfulnessDenominator']
    df.isnull().sum()
    df['Helpfulness'] = df['Helpfulness'].fillna(0)

    df['ReviewLength'] = df['Text'].apply(lambda x: len(x.split()) if isinstance(x,
```

```

df['ReviewCharLength'] = df['Text'].apply(lambda x: len(x) if isinstance(x, str) else 0)

# Count of capitalized words
df['CapitalizedCount'] = df['Text'].apply(lambda x: sum(map(str.isupper, x.split())))

# Sentiment Analysis
df['Polarity'] = df['Text'].apply(lambda x: TextBlob(x).sentiment.polarity if isinstance(x, str) else 0)
df['Subjectivity'] = df['Text'].apply(lambda x: TextBlob(x).sentiment.subjectivity if isinstance(x, str) else 0)

product_encoder = OneHotEncoder()
product_encoded = product_encoder.fit_transform(df[['ProductId']]).toarray()
df = df.join(pd.DataFrame(product_encoded, columns=product_encoder.get_feature_names_out()))

return df

# Load the dataset
trainingSet = pd.read_csv("../data/train.csv")
#####
#分子分母?
trainingSet['Helpfulness'] = trainingSet['HelpfulnessNumerator'] / trainingSet['HelpfulnessDenominator']
mean_denominator = trainingSet[trainingSet['Helpfulness'] != 0]['Helpfulness'].mean()
trainingSet['Helpfulness'] = trainingSet['Helpfulness'].fillna(mean_denominator)
trainingSet['ReviewLength'] = trainingSet.apply(lambda row: len(row['Text'].split()) if isinstance(row['Text'], str) else 0)
#####

#####
#'ProductId'
#LabelEncoder
label_encoder = LabelEncoder()
#ProductId to code
trainingSet['ProductId_encoded'] = label_encoder.fit_transform(trainingSet['ProductId'])
#####

#####
good_words = ['interesting', 'Interesting', 'classic', 'Classic', 'best', 'Best', 'enjoyable', 'Enjoyable', 'great', 'Great', 'amazing', 'Amazing', 'fantastic', 'Fantastic', 'awesome', 'Awesome', 'wonderful', 'Wonderful', 'excellent', 'Excellent', 'superb', 'Superb', 'brilliant', 'Brilliant', 'outstanding', 'Outstanding', 'phenomenal', 'Phenomenal', 'incredible', 'Incredible', 'unbelievable', 'Unbelievable', 'amazing', 'Amazing', 'fantastic', 'Fantastic', 'awesome', 'Awesome', 'wonderful', 'Wonderful', 'excellent', 'Excellent', 'superb', 'Superb', 'brilliant', 'Brilliant', 'outstanding', 'Outstanding', 'phenomenal', 'Phenomenal', 'incredible', 'Incredible', 'unbelievable', 'Unbelievable']
bad_words = ['superficial', 'Superficial', 'terrible', 'Terrible', 'obnoxious', 'Obnoxious', 'boring', 'Boring', 'dull', 'Dull', 'waste of time', 'Waste of time', 'disappointing', 'Disappointing', 'mediocre', 'Mediocre', 'forgettable', 'Forgettable', 'predictable', 'Predictable', 'cliché', 'Cliché', 'overrated', 'Overrated', 'underwhelming', 'Underwhelming', 'lackluster', 'Lackluster', 'forgettable', 'Forgettable', 'predictable', 'Predictable', 'cliché', 'Cliché', 'overrated', 'Overrated', 'underwhelming', 'Underwhelming', 'lackluster', 'Lackluster']
negation = ['not', 'Not', 'too little', 'Too little', "couldn't", "Couldn't", "can't", "Can't", "won't", "Won't", "isn't", "Isn't", "aren't", "Aren't", "wasn't", "Wasn't", "weren't", "Weren't", "shouldn't", "Shouldn't", "mustn't", "Mustn't", "oughtn't", "Oughtn't", "mayn't", "Mayn't", "mightn't", "Mightn't", "shan't", "Shan't", "won't", "Won't", "isn't", "Isn't", "aren't", "Aren't", "wasn't", "Wasn't", "weren't", "Weren't", "shouldn't", "Shouldn't", "mustn't", "Mustn't", "oughtn't", "Oughtn't", "mayn't", "Mayn't", "mightn't", "Mightn't", "shan't", "Shan't"]

trainingSet["negation_count"] = trainingSet['Text'].str.count('|'.join(negation))
trainingSet['negation_count'] = trainingSet['negation_count'].fillna(0)
trainingSet["good_count"] = trainingSet['Text'].str.count('|'.join(good_words))
trainingSet['good_count'] = trainingSet['good_count'].fillna(0)
trainingSet["bad_count"] = trainingSet['Text'].str.count('|'.join(bad_words))
trainingSet['bad_count'] = trainingSet['bad_count'].fillna(0)
#####

#####
trainingSet['same_user'] = trainingSet.duplicated(subset=['UserId'])
trainingSet['returning_user'] = trainingSet.apply(lambda row: 1 if row['same_user'] else 0)
trainingSet['same_movie'] = trainingSet.duplicated(subset=['ProductId'])
trainingSet['movie_repeated'] = trainingSet.apply(lambda row: 1 if row['same_movie'] else 0)
#####

```

```
#####
# Review length in characters
trainingSet['ReviewCharLength'] = trainingSet['Text'].apply(lambda x: len(x) if isi
# Count of capitalized words
trainingSet['CapitalizedCount'] = trainingSet['Text'].apply(lambda x: sum(map(str.i
# Sentiment Analysis
trainingSet['Polarity'] = trainingSet['Text'].apply(lambda x: TextBlob(x).sentiment
trainingSet['Subjectivity'] = trainingSet['Text'].apply(lambda x: TextBlob(x).senti
#####

#####
#score取1/5
#####

# Process the DataFrame
train_processed = trainingSet

#####
# n_5=train_processed[train_processed['Score']==5].shape[0]//2
# df_5=train_processed[train_processed['Score']==5].sample(n=n_5,random_state=1)

# df_not5=train_processed[train_processed['Score']!=5]
# train_processed_balenced=pd.concat([df_5,df_not5])
# train_processed_balenced=train_processed_balenced.sample(frac=1,random_state=1).res
#####

# Load test set
submissionSet = pd.read_csv("./data/test.csv")

# Merge on Id so that the test set can have feature columns as well
#testX= pd.merge(train_processed_balenced, submissionSet, left_on='Id', right_on='Id')
testX= pd.merge(train_processed, submissionSet, left_on='Id', right_on='Id')

testX = testX.drop(columns=['Score_x'])
testX = testX.rename(columns={'Score_y': 'Score'})

# The training set is where the score is not null
trainX = train_processed[train_processed['Score'].notnull()]

# Save the datasets with the new features for easy access later
testX.to_csv("./data/X_test.csv", index=False)
trainX.to_csv("./data/X_train.csv", index=False)
```

```
In [ ]: trainX.shape
```

```
Out[ ]: (122283, 23)
```

Creating your model

```
In [ ]: import pickle
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn.neighbors import KNeighborsClassifier
```

```

from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score, confusion_matrix, mean_squared_error
from sklearn.linear_model import LinearRegression
from sklearn.linear_model import Ridge, HuberRegressor
from sklearn.ensemble import RandomForestRegressor, GradientBoostingRegressor
from lightgbm import LGBMClassifier

# Load training set with new features into DataFrame
X_train = pd.read_csv("./data/X_train.csv")

# Split training set into training and testing set
X_train, X_test, Y_train, Y_test = train_test_split(
    X_train.drop(['Score'], axis=1),
    X_train['Score'],
    test_size=1/4.0,
    random_state=0
)

# This is where you can do more feature selection
X_train_processed = X_train.drop(columns=['Id', 'ProductId', 'UserId', 'Text', 'Summary'])
X_test_processed = X_test.drop(columns=['Id', 'ProductId', 'UserId', 'Text', 'Summary'])

#bigger number means faster
#model=RandomForestRegressor()#2
model=GradientBoostingRegressor()#3
#model = LinearRegression()#1
#model=Ridge(alpha=7)#2
#model=HuberRegressor(alpha=2)#1

model.fit(X_train_processed, Y_train)

# Save the model
with open('linear_regression_model.pkl', 'wb') as f:
    pickle.dump(model, f)

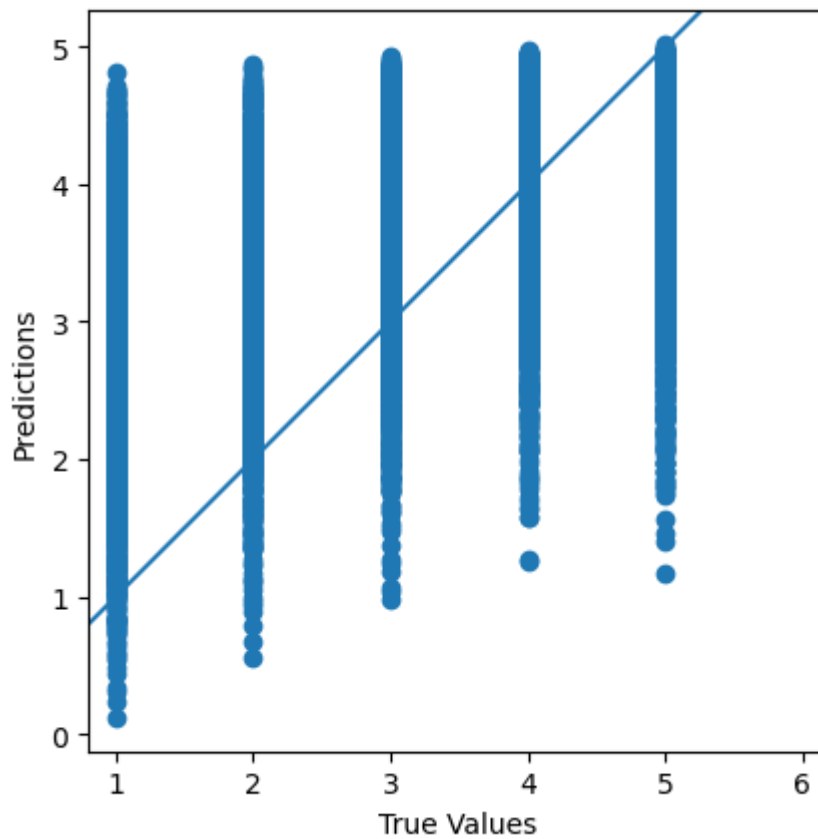
# Evaluate your model on the testing set
Y_test_predictions = model.predict(X_test_processed)

print("RMSE on testing set = ", mean_squared_error(Y_test, Y_test_predictions)**(1/2))

# Since it's a regression problem, you can plot the actual vs predicted values
plt.scatter(Y_test, Y_test_predictions)
plt.xlabel('True Values')
plt.ylabel('Predictions')
plt.axis('equal')
plt.axis('square')
plt.plot([-100, 100], [-100, 100])
plt.show()

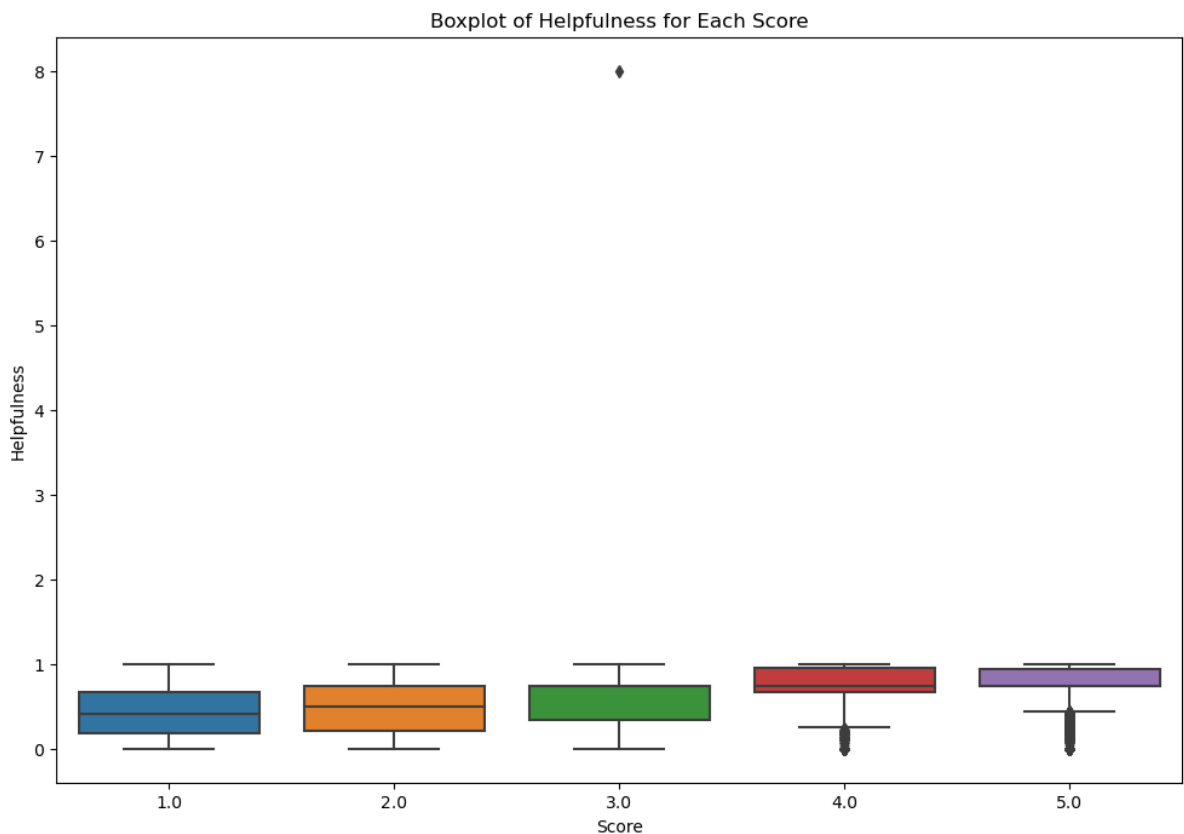
```

RMSE on testing set = 0.9322875053154894



```
In [ ]: # plt.figure(figsize=(10, 8))
# plt.hexbin(trainX['Helpfulness'], trainX['Score'], gridsize=50, cmap='Blues')
# plt.colorbar(label='Count in bin')
# plt.xlabel('Helpfulness')
# plt.ylabel('Score')
# plt.title('Hexbin Plot of Helpfulness vs Score')
# plt.show()
```

```
In [ ]: plt.figure(figsize=(12, 8))
sns.boxplot(x='Score', y='Helpfulness', data=trainX)
plt.title('Boxplot of Helpfulness for Each Score')
plt.xlabel('Score')
plt.ylabel('Helpfulness')
plt.show()
```



```
In [ ]: testX.isnull().sum()
```

```
Out[ ]: Id                                0
ProductId                               0
UserId                                  0
HelpfulnessNumerator                     0
HelpfulnessDenominator                   0
Time                                     0
Summary                                 0
Text                                    0
Helpfulness                             0
ReviewLength                            0
ProductId_encoded                        0
negation_count                           0
good_count                              0
bad_count                                0
same_user                               0
returning_user                           0
same_movie                              0
movie_repeated                           0
ReviewCharLength                         0
CapitalizedCount                         0
Polarity                                 0
Subjectivity                             0
Score                                   17470
dtype: int64
```

Create the Kaggle submission

```
In [ ]: X_submission = pd.read_csv("../data/X_test.csv")
X_submission_processed = X_submission.drop(columns=['Id', 'ProductId', 'UserId', 'Te

X_submission['Score'] = model.predict(X_submission_processed)
submission = X_submission[['Id', 'Score']]
submission.to_csv("../data/submission.csv", index=False)
```


Now you can upload the `submission.csv` to kaggle