# Exploration

```python
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, alpha=.5)
plt.title("Top 25 most rated Products")
plt.show()

trainingSet['ProductId'].value_counts().nsmallest(25).plot(kind='bar',
legend=True, alpha=.5)
plt.title("Top 25 least rated Products")
plt.show()

trainingSet['UserId'].value_counts().nlargest(25).plot(kind='bar',
legend=True, alpha=.5)
plt.title("Top 25 Reviewers")
plt.show()

trainingSet['UserId'].value_counts().nsmallest(25).plot(kind='bar',
legend=True, alpha=.5)
plt.title("Lowest 25 Reviewers")
plt.show()

trainingSet[['Score',
'HelpfulnessNumerator']].groupby('Score').mean().plot(kind='bar',
legend=True, alpha=.5)
plt.title("Mean Helpfulness Numerator per Score")
```

```
plt.show()

trainingSet[['Score',
'ProductId']].groupby('ProductId').mean().nlargest(25,
'Score').plot(kind='bar', legend=True, alpha=.5)
plt.title("Top 25 best rated Products")
plt.show()

trainingSet[['Score',
'ProductId']].groupby('ProductId').mean().nsmallest(25,
'Score').plot(kind='bar', legend=True, alpha=.5)
plt.title("Top 25 worst rated Products")
plt.show()

trainingSet[['Score', 'UserId']].groupby('UserId').mean().nlargest(25,
'Score').plot(kind='bar', legend=True, alpha=.5)
plt.title("Top 25 kindest Reviewers")
plt.show()

trainingSet[['Score',
'UserId']].groupby('UserId').mean().nsmallest(25,
'Score').plot(kind='bar', legend=True, alpha=.5)
plt.title("Top 25 harshest Reviewers")
plt.show()

trainingSet[trainingSet['ProductId'].isin(trainingSet['ProductId'].val
ue_counts().nlargest(25).index.tolist())][['Score',
'ProductId']].groupby('ProductId').mean().plot(kind='bar',
legend=True, alpha=.5)
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 | A3VLX5Z090RQ0V | 1 |
| 1 | 1632470 | B00BEIYSL4 | AUDXDMFM49NGY | 0 |
| 2 | 9771 | 0767809335 | A3LFIA97BUU5IE | 3 |
| 3 | 218855 | 6300215792 | A1QZM75342ZQVQ | 1 |
| 4 | 936225 | B000B5XOZW | 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 LAAAAADEEE!!!!
4   Herzog the Great Traveler of both natural and ...
```
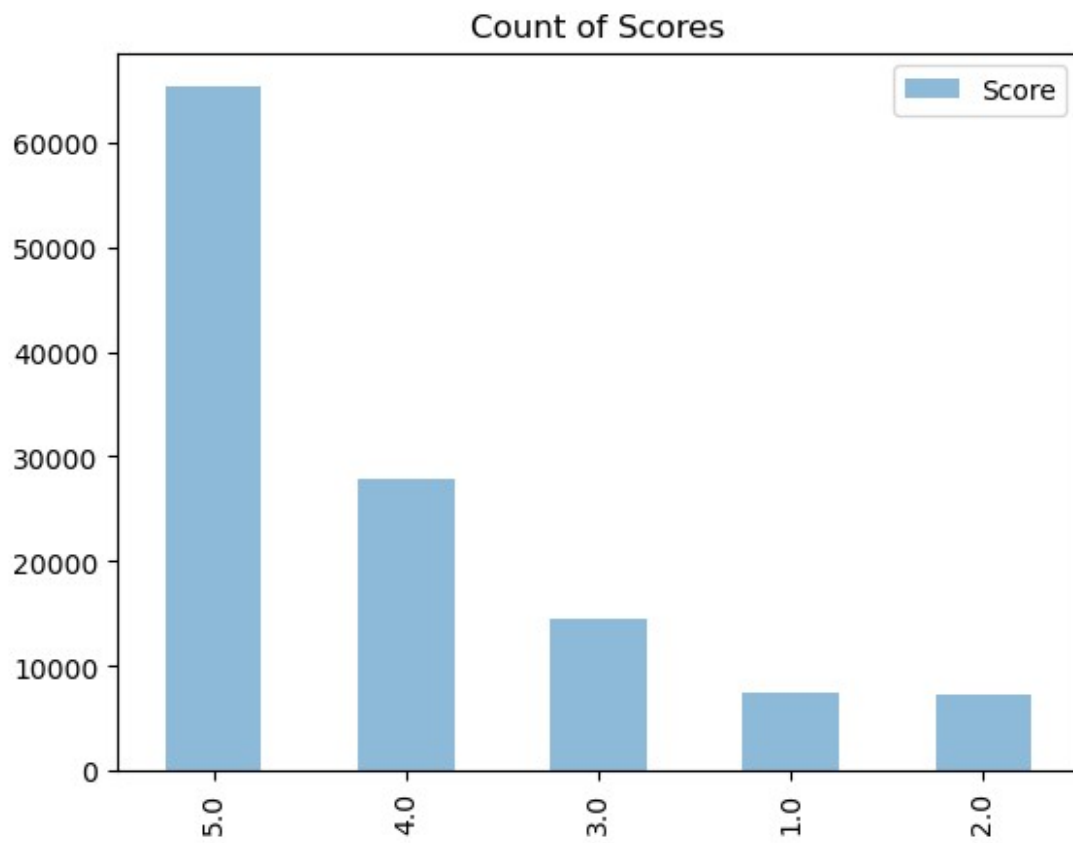
```
                                                 Text  Score
0   I was very anxious to see the Uncut version of...    2.0
1                          Movie was okay...not great.    3.0
2   Like the Dinosaur Collector's Edition DVD, thi...    1.0
3   Come on, now.....  this has to be, by far, the...    5.0
4   I've always been a great admirer of Herzog's o...    4.0
```
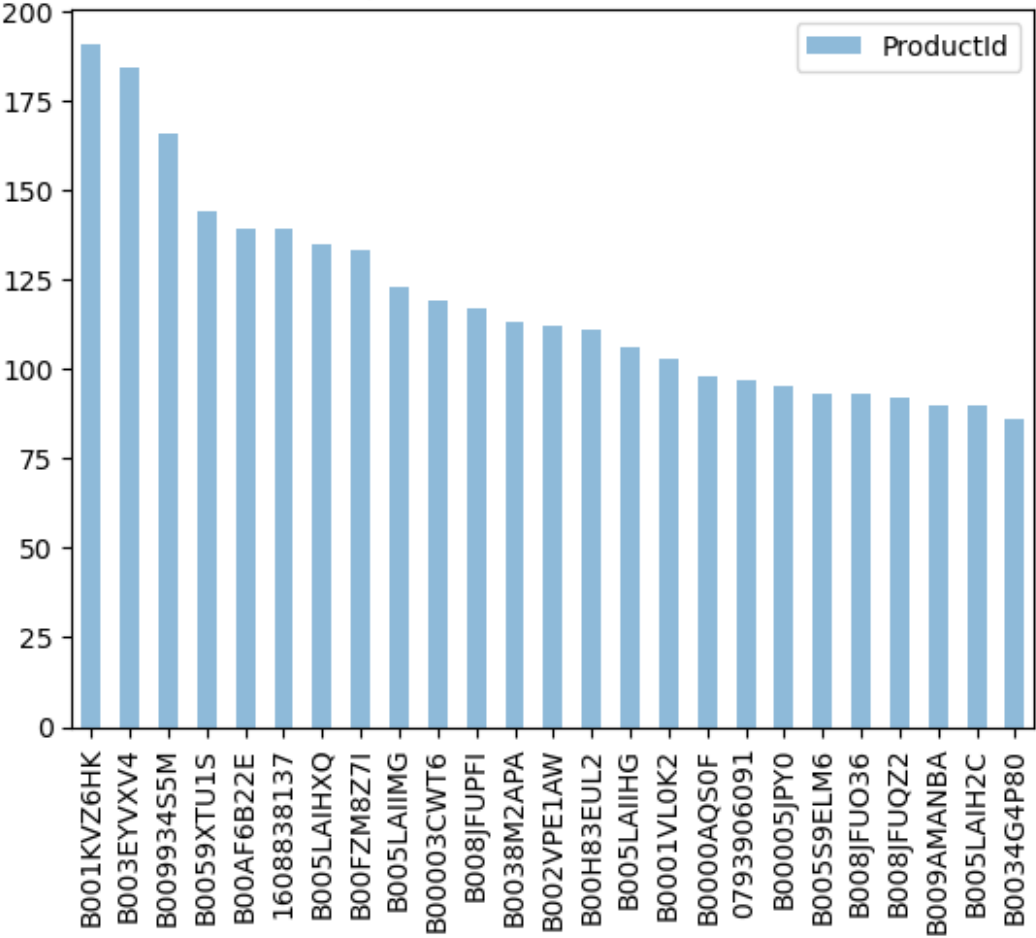
```
         Id  Score
0    786781    NaN
1     17153    NaN
2   1557328    NaN
3   1242666    NaN
4   1359242    NaN
```

```
                 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
```
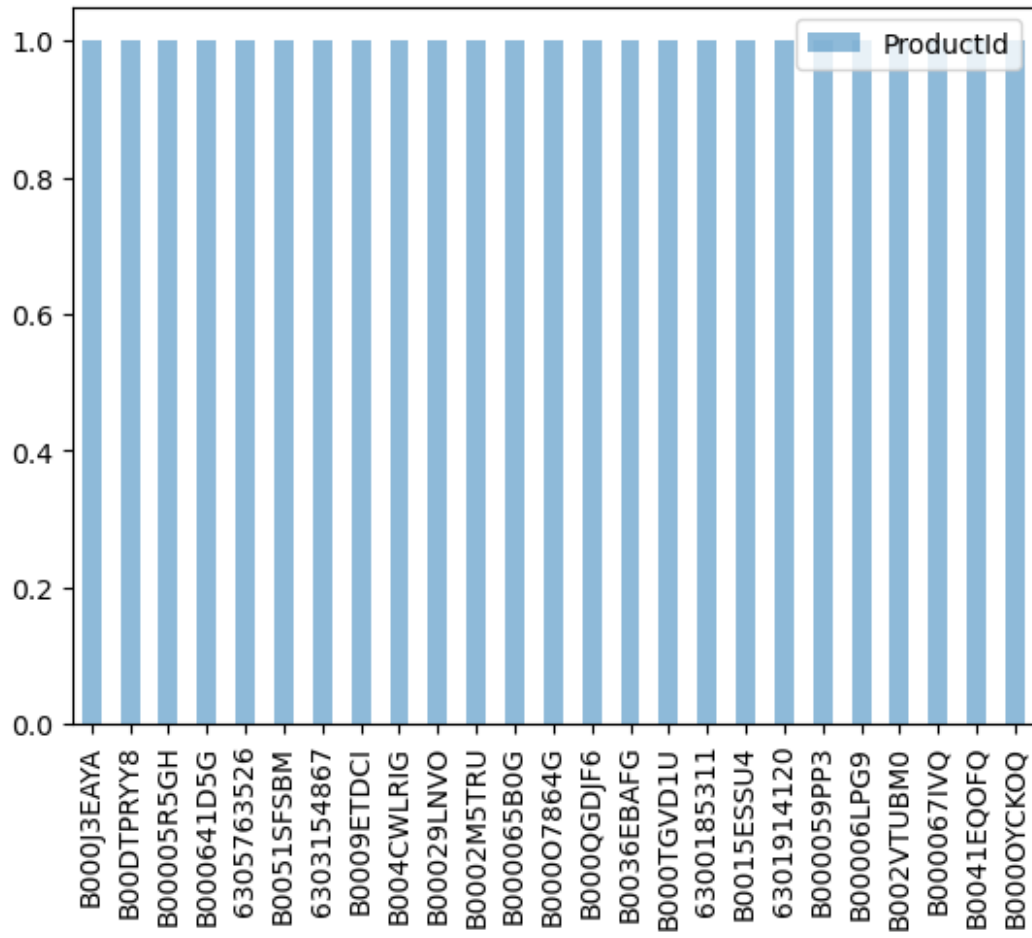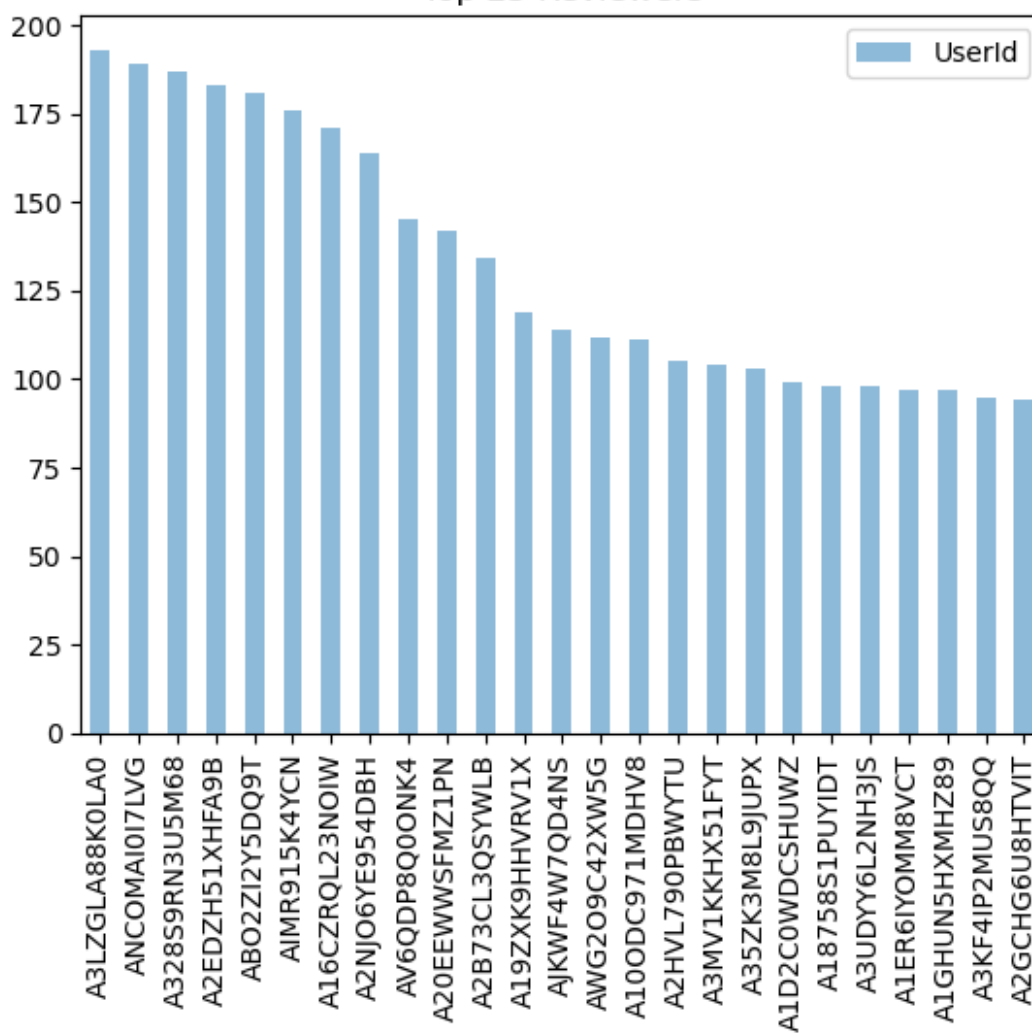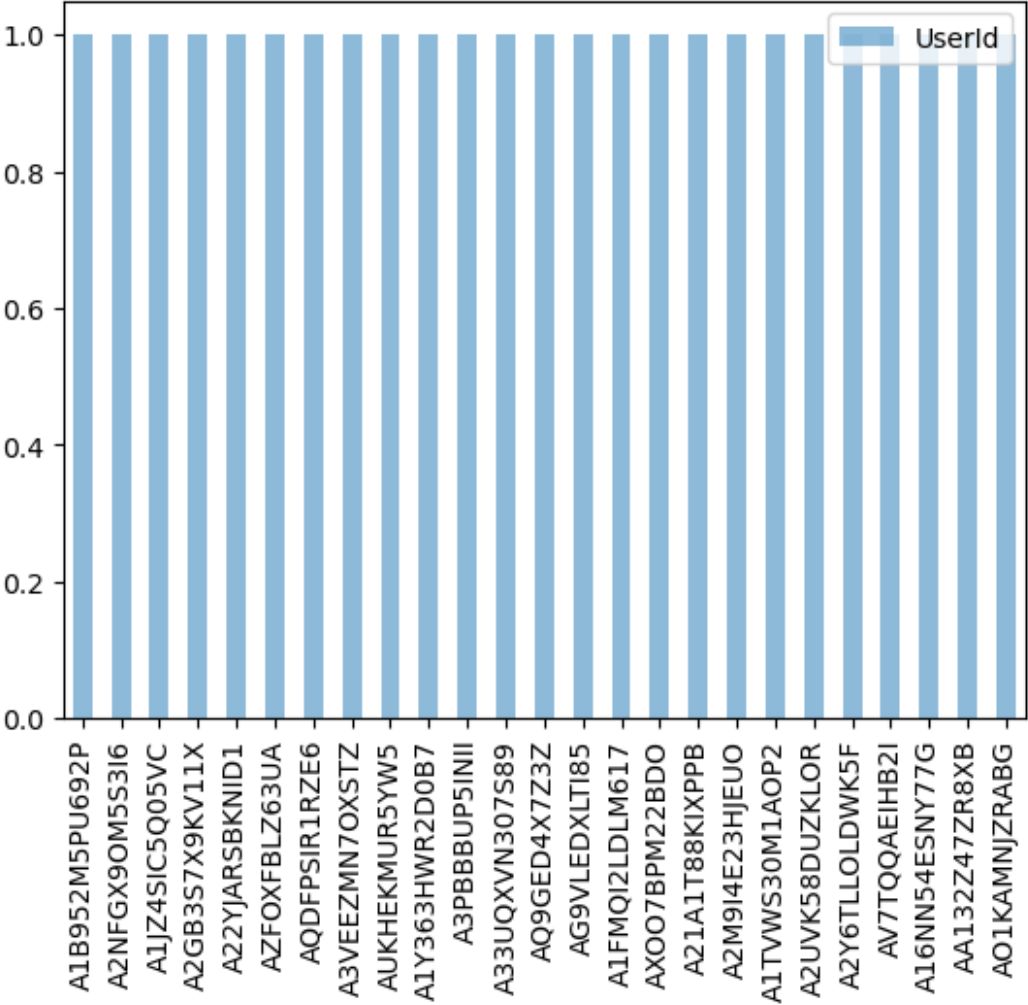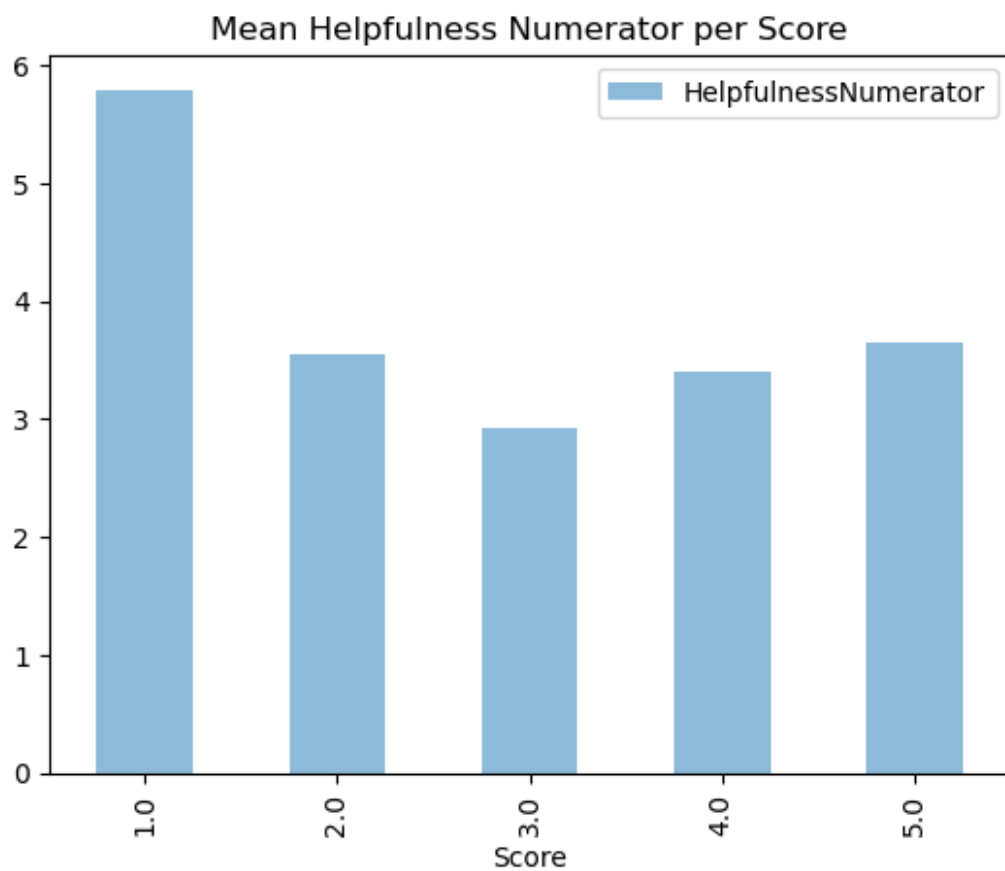
Count of Scores

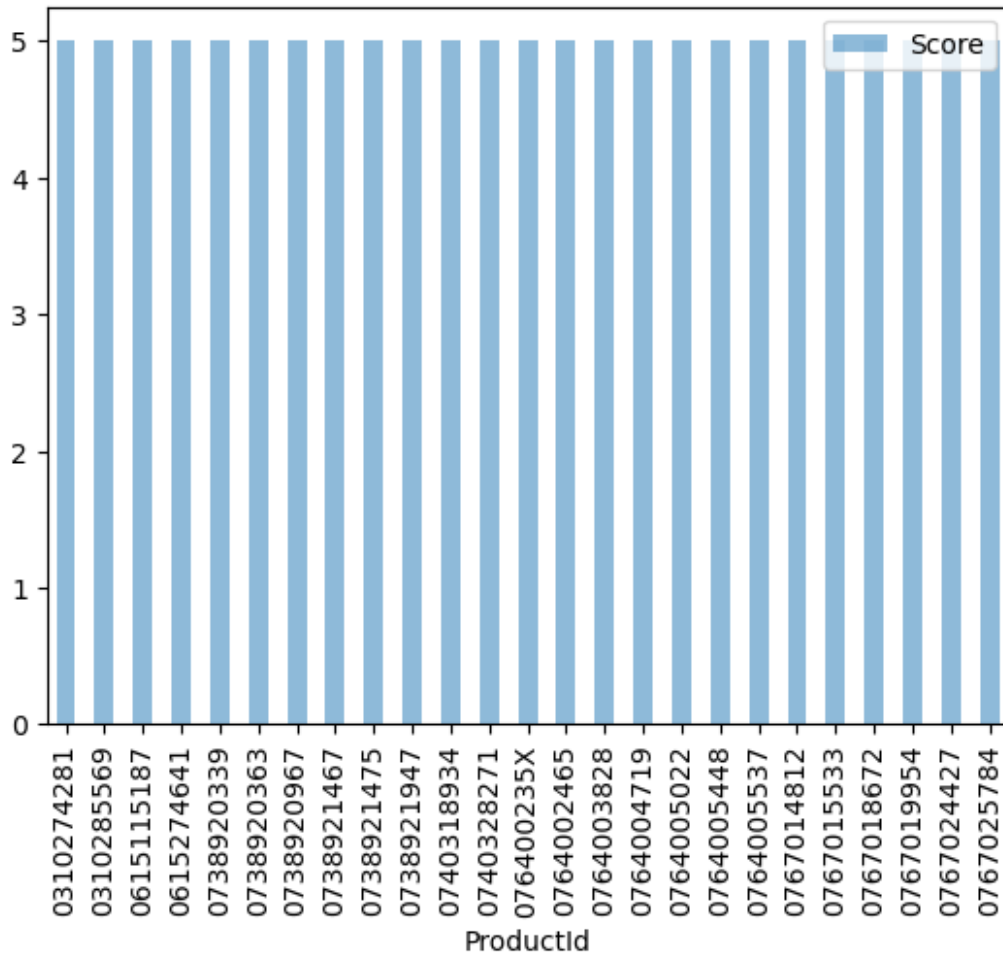Top 25 most rated Products

Top 25 least rated Products

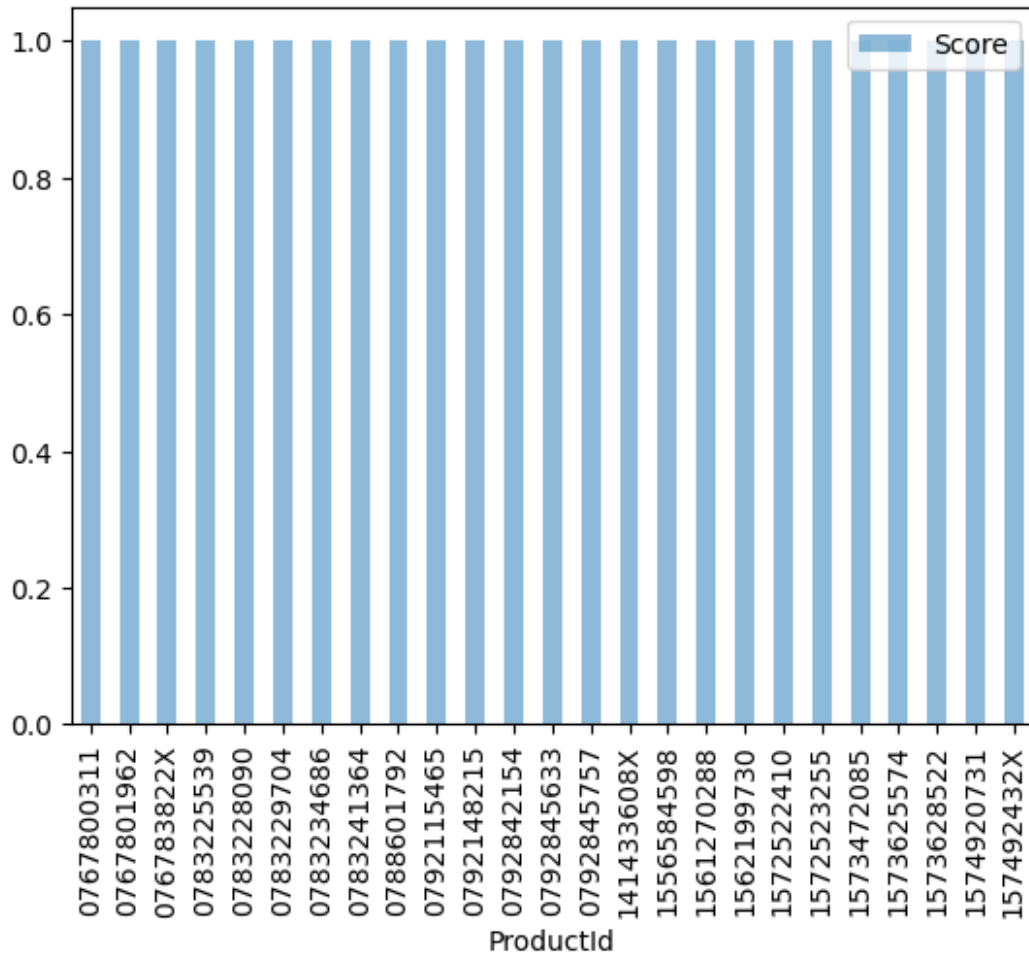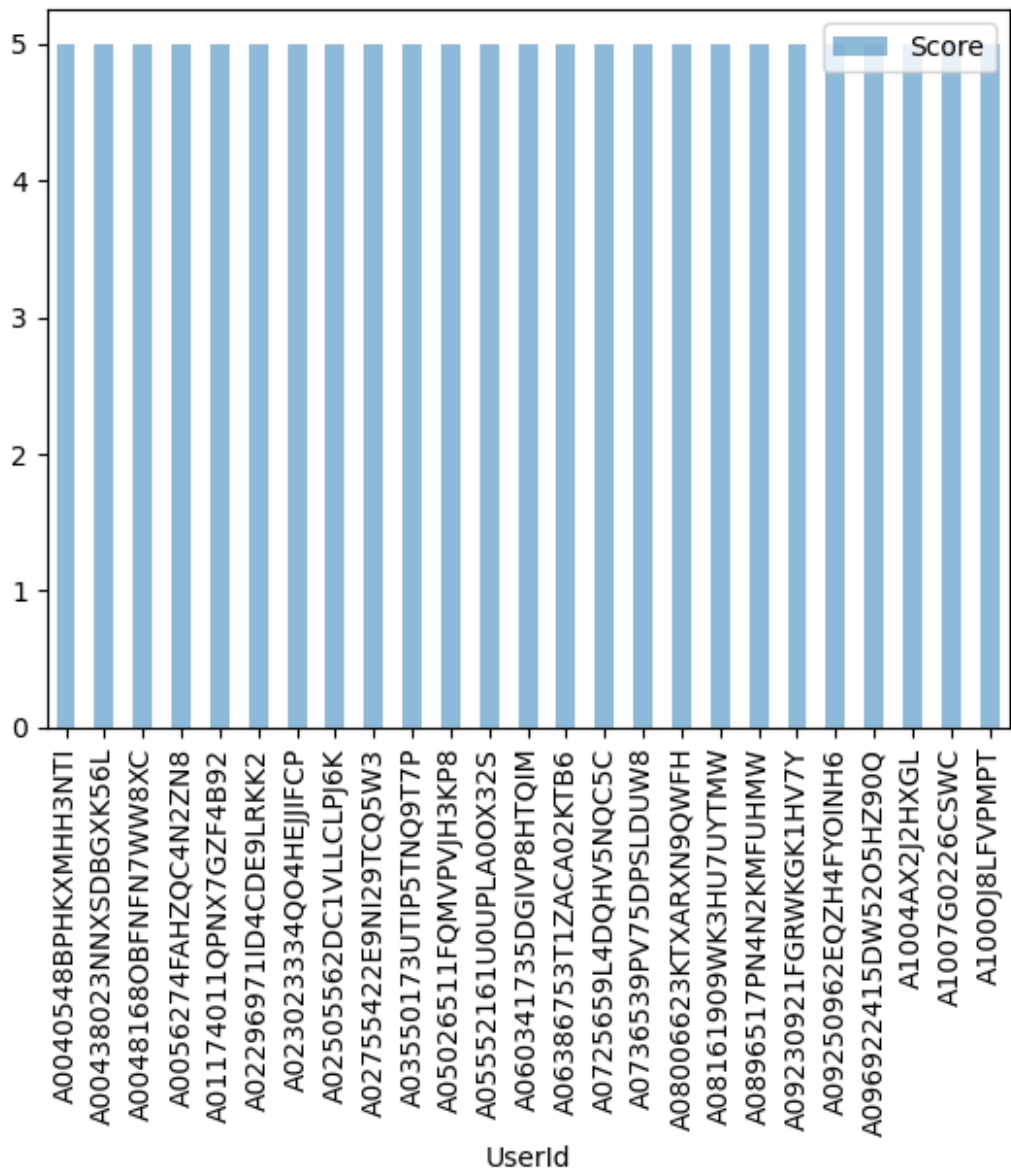Top 25 Reviewers

Lowest 25 Reviewers

Mean Helpfulness Numerator per Score

## Top 25 best rated Products



| ProductId |

# Top 25 worst rated Products



Score

| ProductId | |

Top 25 kindest Reviewers

Top 25 harshest Reviewers

Mean of top 25 most rated Products

## Feature Extraction

```python
import pandas as pd

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

    df['Helpfulness'] = df['HelpfulnessNumerator'] /
df['HelpfulnessDenominator']
    df['Helpfulness'] = df['Helpfulness'].fillna(0)

    df['ReviewLength'] = df.apply(lambda row :
len(row['Text'].split()) if type(row['Text']) == str else 0, axis = 1)

    return df

# Load the dataset
trainingSet = pd.read_csv("./data/train.csv")
```

```
# Process the DataFrame
train_processed = process(trainingSet)

# 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, 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)
```

## Creating your model

```
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

# 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=3
    )

# 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'])

# Learn the model
model = KNeighborsClassifier(n_neighbors=20).fit(X_train_processed,
Y_train)
```
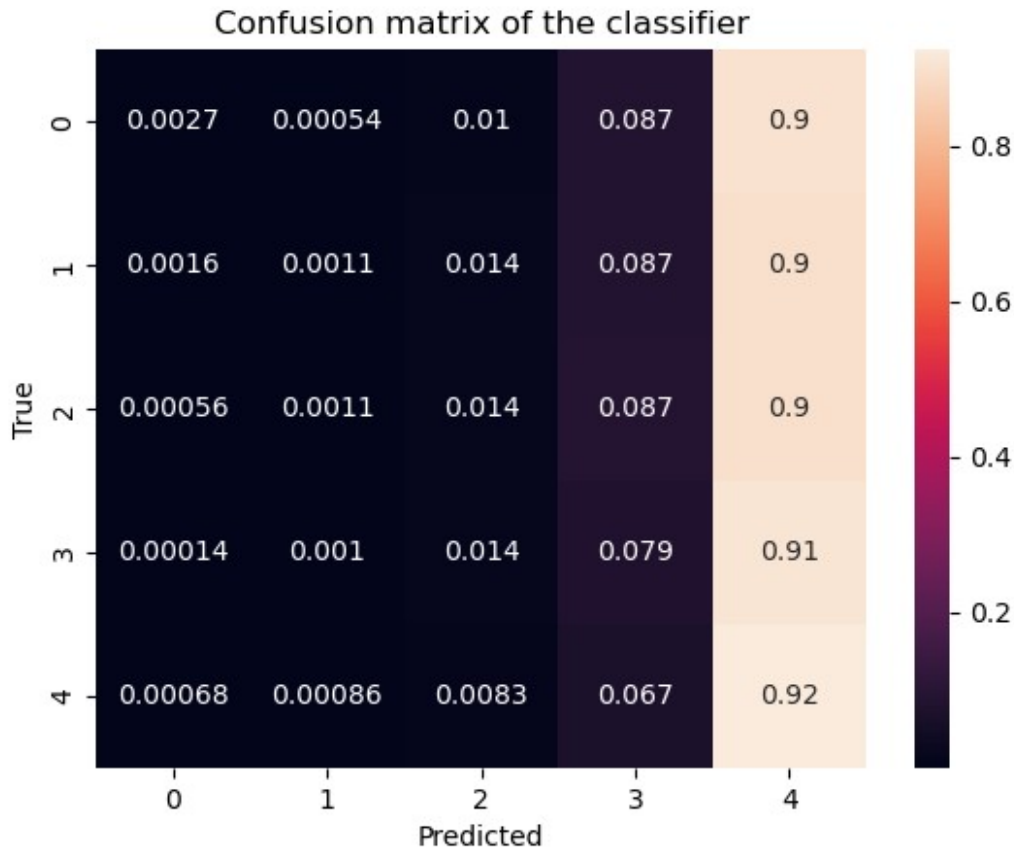
```python
# pickle model - saves it so you can load it later
with open('knn_20_model.obj', 'wb') as f:
        pickle.dump(model, f)
# to load pickled model:
# with open('filename', 'rb') as f:
#     model = pickle.load(f)

# Evaluate your model on the testing set
Y_test_predictions = model.predict(X_test_processed)
print("Accuracy on testing set = ", accuracy_score(Y_test,
Y_test_predictions))
print("RMSE on testing set = ", mean_squared_error(Y_test,
Y_test_predictions) ** 0.5)

# Plot a confusion matrix
cm = confusion_matrix(Y_test, Y_test_predictions, normalize='true')
sns.heatmap(cm, annot=True)
plt.title('Confusion matrix of the classifier')
plt.xlabel('Predicted')
plt.ylabel('True')
plt.show()

Accuracy on testing set =  0.5121520395145727
RMSE on testing set =  1.4600992087486073
```

Confusion matrix of the classifier

## Create the Kaggle submission

```
X_submission = pd.read_csv("./data/X_test.csv")
X_submission_processed = X_submission.drop(columns=['Id', 'ProductId',
'UserId', 'Text', 'Summary', 'Score'])

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