

Context: Online E-commerce websites like Amazon, Flipkart uses different recommendation models to provide different suggestions to different users. Amazon currently uses item-to-item collaborative filtering, which scales to massive data sets and produces high-quality recommendations in real-time.

Attribute Information:

- userId : Every user identified with a unique id
- productId : Every product identified with a unique id
- Rating : Rating of the corresponding product by the corresponding user
- timestamp : Time of the rating

Objective: Build a recommendation system to recommend products to customers based on the their previous ratings for other products.

1. Read and explore the given dataset. (Rename column/add headers, plot histograms, find data characteristics)

Import Libraries

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split

%matplotlib inline
```

Load the Dataset and Add headers

```
In [2]: electronics_data = pd.read_csv("ratings_Electronics.csv")
electronics_data.columns = ['userId', 'productId', 'Rating', 'timestamp']
electronics_data.head()
```

```
Out[2]:      userId  productId  Rating  timestamp
0    AZCOTLJUH82NDG  321732944      5  1341108800
1    AZNWSAGRHCP8N5  439886341      1  1367193600
2    AZWNBD03WNDNKT  439886341      3  1374451200
3    A1GI0U4ZRIABWN  439886341      1  1334707200
4    A1QGNMCG01VW39  511189877      5  1397433600
```

```
In [3]: # Display the data
electronics_data.head()
```

```
Out[3]:      userId  productId  Rating  timestamp
0    AZCOTLJUH82NDG  321732944      5  1341108800
1    AZNWSAGRHCP8N5  439886341      1  1367193600
2    AZWNBD03WNDNKT  439886341      3  1374451200
3    A1GI0U4ZRIABWN  439886341      1  1334707200
4    A1QGNMCG01VW39  511189877      5  1397433600
```

```
In [4]: #Shape of the data
electronics_data.shape
```

```
Out[4]: (1048575, 4)
```

```
In [5]: #Faking subset of the dataset
electronics_data=electronics_data.iloc[:1048576,0:]
```

```
In [6]: #Check the datatypes
electronics_data.dtypes
```

```
Out[6]: userId      object
productId  object
Rating      int64
timestamp  int64
dtype: object
```

```
In [7]: electronics_data.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1048575 entries, 0 to 1048574
Data columns (total 4 columns):
 #   Column      Non-Null Count  Dtype
---  ---
 0   userId     1048575 non-null  object
 1   productId  1048575 non-null  object
 2   Rating     1048575 non-null  int64
 3   timestamp  1048575 non-null  int64
dtypes: int64(2), object(2)
memory usage: 32.0+ MB
```

```
In [8]: #Five point summary
electronics_data.describe()['Rating'].T
```

```
Out[8]: count      1.048575e+06
mean       3.973379e+00
std        1.399329e+00
min        1.000000e+00
25%        3.000000e+00
50%        5.000000e+00
75%        5.000000e+00
max        5.000000e+00
Name: Rating, dtype: float64
```

```
In [9]: #Find the minimum and maximum ratings
print('Minimum rating is: %d' % (electronics_data.Rating.min()))
print('Maximum rating is: %d' % (electronics_data.Rating.max()))

Minimum rating is: 1
Maximum rating is: 5

The rating of the product range from 0 to 1
```

Handling Missing values

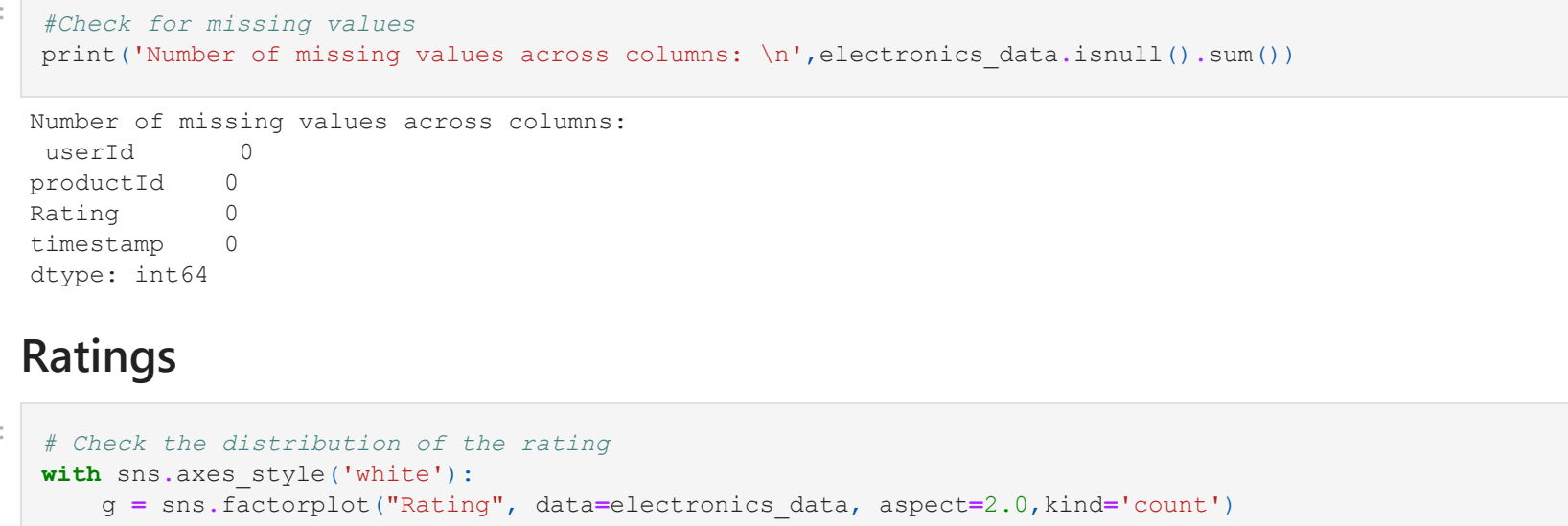
```
In [10]: #Check for missing values
print('Number of missing values across columns: \n',electronics_data.isnull().sum())

Number of missing values across columns:
userId      0
productId   0
Rating      0
timestamp   0
dtype: int64
```

Ratings

```
In [11]: # Check the distribution of the rating
with sns.axes_style('white'):
    g = sns.FacetGrid("Rating", data=electronics_data, aspect=2.0, kind='count')
    g.set_ylabels('Total number of ratings')
```

C:\Users\Prashant Varshney\anaconda3\lib\site-packages\seaborn\categorical.py:3717: UserWarning: The 'factorplot' function has been renamed to 'catplot'. The original name will be removed in a future release. Please update your code. Note that the default 'kind' in 'factorplot' ('point') has changed 'strip' in 'catplot'.
Warnings: Warn(nmsg)
C:\Users\Prashant Varshney\anaconda3\lib\site-packages\seaborn_decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be 'data', and passing other arguments without an explicit keyword will result in an error or misinterpretation.
Warnings: Warn(n



Most of the people has given the rating of 5

Unique Users and products

```
In [12]: print("Total data ")
print("n="*50)
print("Total no of ratings :",electronics_data.shape[0])
print("Total No of Users  :", len(np.unique(electronics_data.userId)))
print("Total No of products :", len(np.unique(electronics_data.productId)))

Total data
-----
```

```
Total no of ratings : 1048575
Total No of Users   : 786329
Total No of products : 61893
```

Dropping the TimeStamp Column

```
In [13]: #dropping the Timestamp column
electronics_data.drop(['timestamp'], axis=1,inplace=True)
```

Analyzing the rating

```
In [14]: #Analysis of rating given by the user
no_of_rated_products_per_user = electronics_data.groupby(by='userId')['Rating'].count().sort_values(ascending=False)
no_of_rated_products_per_user.head()
```

```
Out[14]: userId      412
A5JLAUZARJ0BO      146
A231NM22JL0U3      249
A2580V8868EA       144
A6FLAB28I879       146
AT6C8DCE4TRGA      128
Name: Rating, dtype: int64
```

```
In [15]: no_of_rated_products_per_user.describe()
```

```
Out[15]: count      786329.000000
mean       1.333507
std        1.385613
min        1.000000
25%        1.000000
50%        1.000000
75%        1.000000
max        412.000000
Name: Rating, dtype: float64
```

```
In [16]: quantiles = no_of_rated_products_per_user.quantile(np.arange(0,1.01,0.01), interpolation='higher')
```

```
In [17]: plt.figure(figsize=(10,10))
plt.title("Quantiles and their Values")
quantiles.plot()
# quantiles with 0.05 difference
plt.scatter(x=quantiles.index[:5], y=quantiles.values[:5], c='orange', label='quantiles with 0.05 intervals')
# quantiles with 0.25 difference
plt.scatter(x=quantiles.index[:25], y=quantiles.values[:25], c='m', label = "quantiles with 0.25 intervals")
plt.xlabel('No of ratings by user')
plt.ylabel('Value at the quantile')
plt.legend(loc='best')
plt.show()
```



2. Take a subset of the dataset to make it less sparse/ denser. (For example, keep the users only who has given 50 or more number of ratings)

```
In [18]: print("\n No of rated product more than 50 per user : ( )\n".format(sum(no_of_rated_products_per_user >= 50) ) )

No of rated product more than 50 per user : 38
```

```
In [19]: #getting the new dataframe which contains users who has given 50 or more ratings
new_df=electronics_data.groupby("productId").filter(lambda x:x['Rating'].count() >=50)
```

3. Split the data randomly into train and test dataset. (For example, split it in 70/30 ratio)

```
In [21]: from surprise import KNNWithMeans
from surprise import Dataset
from surprise import accuracy
from surprise import Reader
import os
from surprise.model_selection import train_test_split
```

```
In [ ]: 
```

```
In [22]: #Reading the dataset
reader = Reader(rating_scale=(1, 5))
data = Dataset.load_from_df(new_df,new_reader)
```

```
In [23]: #Splitting the dataset
trainset, testset = train_test_split(data, test_size=0.3,random_state=10)
```

4. Build Popularity Recommender model.

```
In [24]: no_of_ratings_per_movie = new_df.groupby(by='productId')['Rating'].count().sort_values(ascending=False)
```

```
fig = plt.figure(figsize=plt.figaspect(.5))
ax = plt.gca()
plt.plot(no_of_ratings_per_movie.values)
plt.title("# RATINGS per Product")
plt.xlabel('Product')
plt.ylabel('No of Users who rated a product')
ax.set_xticklabels([])

plt.show()
```



```
In [25]: #Average rating of the product
new_df.groupby('productId')['Rating'].mean().head()
```

```
Out[25]: productId      3.560000
1400501466      4.243902
1400501520      3.884892
1400501776      3.684211
1400526220      3.727273
1400532655      3.727273
Name: Rating, dtype: float64
```

```
In [26]: new_df.groupby('productId')['Rating'].mean().sort_values(ascending=False).head()
```

```
Out[26]: productId      4.947368
B0000DYV98      4.945783
B00009896C      4.885714
B000051L876      4.879310
B0001X3W8      4.869562
Name: Rating, dtype: float64
```

```
In [27]: #Total no of rating for product
new_df.groupby('productId')['Rating'].count().sort_values(ascending=False).head()
```

```
Out[27]: productId      9487
B000013878      5345
B0001PTVEK      4903
B000168B04      4275
B0000787JU      3523
Name: Rating, dtype: int64
```

```
In [28]: ratings_mean_count = pd.DataFrame(new_df.groupby('productId')['Rating'].mean())
```

```
In [29]: ratings_mean_count['rating_counts'] = pd.DataFrame(new_df.groupby('productId')['Rating'].count())
```

```
In [30]: ratings_mean_count.head()
```

```
Out[30]:      Rating  rating_counts
productId
1400501466  3.560000           250
1400501520  4.243902            82
1400501776  3.884892           139
1400526220  3.684211           171
1400532655  3.727273            484
```

```
In [31]: ratings_mean_count['rating_counts'].max()
```

```
Out[31]: 9487
```

```
In [32]: plt.figure(figsize=(8,6))
plt.rcParams['patch.force_edgecolor'] = True
ratings_mean_count['rating_counts'].hist(bins=50)
```

```
Out[32]: <AxesSubplot>
```



```
In [33]: plt.figure(figsize=(8,6))
plt.rcParams['patch.force_edgecolor'] = True
sns.jointplot(x='Rating', y='rating_counts', data=ratings_mean_count, alpha=0.4)
```

```
Out[34]: <seaborn.axisgrid.JointGrid at 0x1f319493a60>
```

```
<Figure size 576x432 with 0 Axes>
```



```
In [35]: popular_products = pd.DataFrame(new_df.groupby('productId')['Rating'].count())
most_popular = popular_products.sort_values('Rating', ascending=False)
most_popular.head(30).plot(kind = "bar")
```

```
Out[35]: <AxesSubplot: xlabel='productId'>
```



5. Build Collaborative Filtering model.

```
In [36]: # Use user_based/true/false to switch between user-based or item-based collaborative filtering
algo = KNNWithMeans(k=5, sim_options={'name': 'pearson_baseline', 'user_based': False})
algo.fit(trainset)
```

Estimating biases using using...
Computing the pearson_baseline similarity matrix...

```
In [36]: Done computing similarity matrix.
<surprise.prediction_algorithms.knns.KNNWithMeans at 0x1f327047100>
```

```
In [37]: # run the trained model against the testset
test_pred = algo.test(testset)
```

```
In [38]: test_pred
```