

RATING PREDICTION

Submitted by
MIENGANDHA SINHA

ACKNOWLEDGMENT

I would like to express my gratitude to the Company FlipRobo Technology to give this project to me. In making of this project I hereby used to take help from the references from some websites and also from which is given by the company as sample documentation and details related project and professional guided me a lot in the project and the other previous projects helped me and guided me with completion of the project.

INTRODUCTION

- Rating prediction

We all know ratings and reviews are important. While they've only been around for about two decades, it's hard to imagine shopping without them. According to consumer research we conducted of 30,000+ global shopping, the majority (**88%**) of shopping use reviews to discover and evaluate products and on based on these the customers easy to find things. The ability to successfully decide whether a review will be helpful to other customers and thus give the product more exposure is vital to companies that support these reviews, e-commercial sites like Amazon, Flipkart, Snapdeal, etc. In this research project, show a new approach to enhance the accuracy of the rating prediction by using machine learning methods the training performance of our model changes as we change the training method, the dataset used for training and the features used in the model. It help to build an application which can predict the rating by seeing the review.

- Conceptual Background of the Domain Problem

A primary objective of this project is to predict the ratings of various products by using some attributes that are highly correlated with a rate and reviews through machine learning models. So they are looking for new machine learning models from new data to help in to predict ratings for the reviews which were written in the past and they don't have a rating. So, we have to build an application which can predict the rating by seeing the review by valuating the model. By doing some research on this project, are able to trained the models and predicting things make the previous background to work efficiently.

● Review of Project

From the collecting data phase to the model building phase gives important variables in the used to predict rating. There all types of reviews and ratings in the datasets gives various information. And then the model building do all the data visualizations, data pre-processing steps, evaluating the model, data cleaning and selecting the best model for the project.

● Motivation for the Problem Undertaken

Here, the datasets have the total of 198 entries with 3 rows, no null values, EDA has to be performed to see whether it gain or loss in the variable and its compare to the price among every aspects, to build machine learning models, to determine the optimal values of Hyper parameters and the selection of the best model, by predicting of the value can help to the clients and for the further change in the market from the new data.

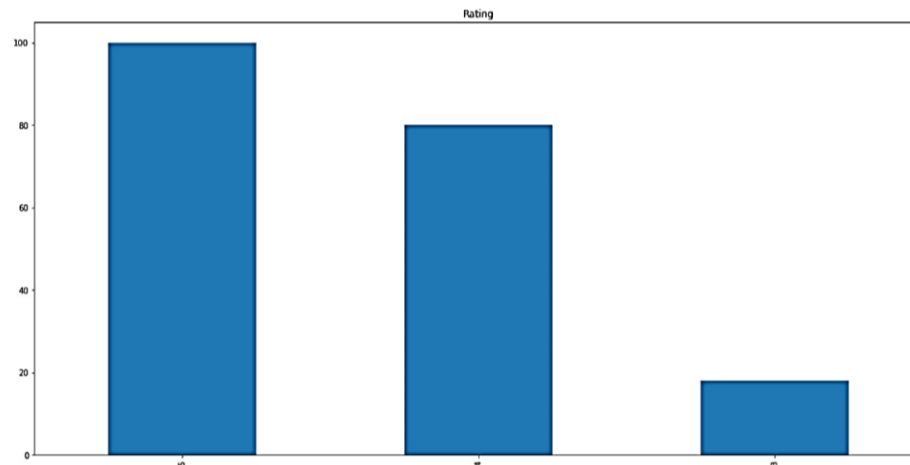
ANALYTICAL PROBLEM FRAMING

- Mathematical/Analytical Modelling of the Problem

In this project, mathematical/analytical modelling are used. Checking the null values found that having no null values in the datasets, the type of data frame is in pandas, data frame info tells that int64(2 variable), object(1 variable), by using the data visualization they are:

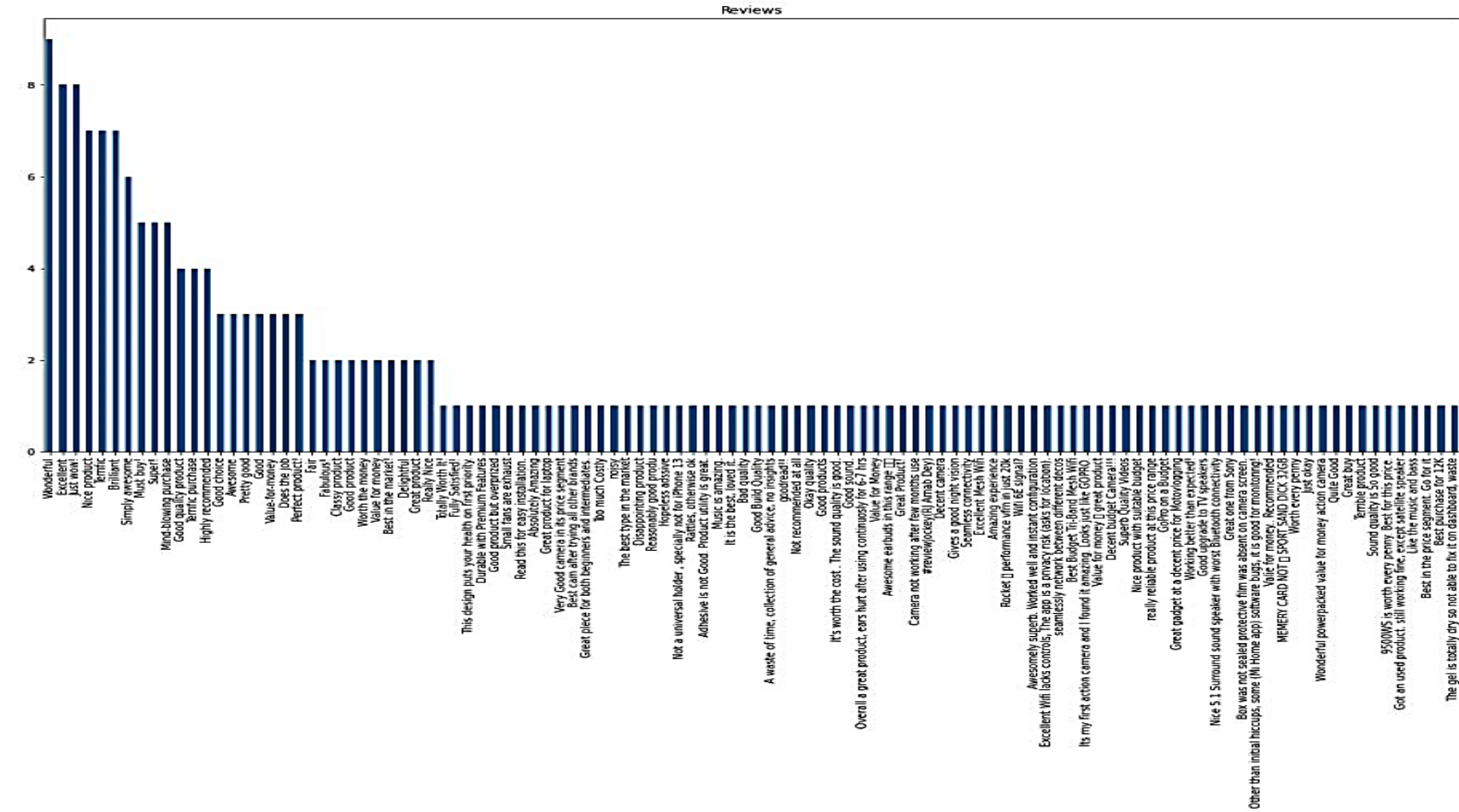
1. It is about the rating of the product:

Five Star rating has the highest accuracy among other stars

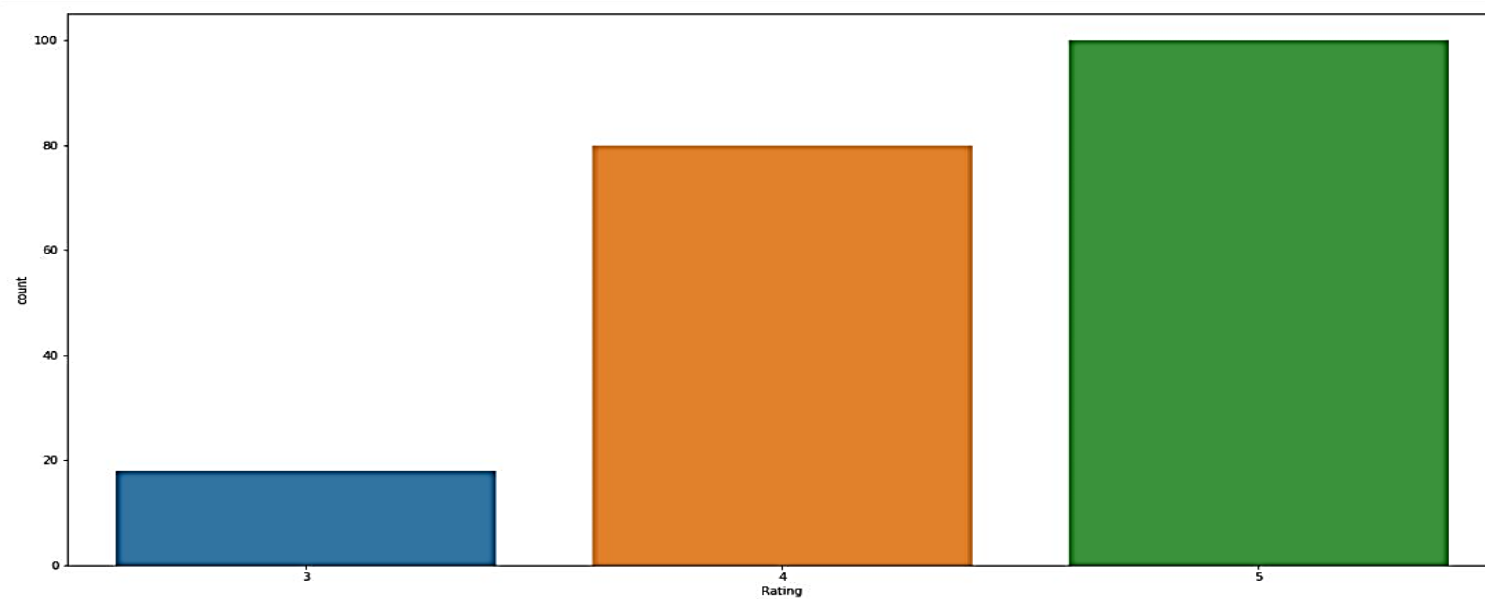


2. It is about the reviews of products:

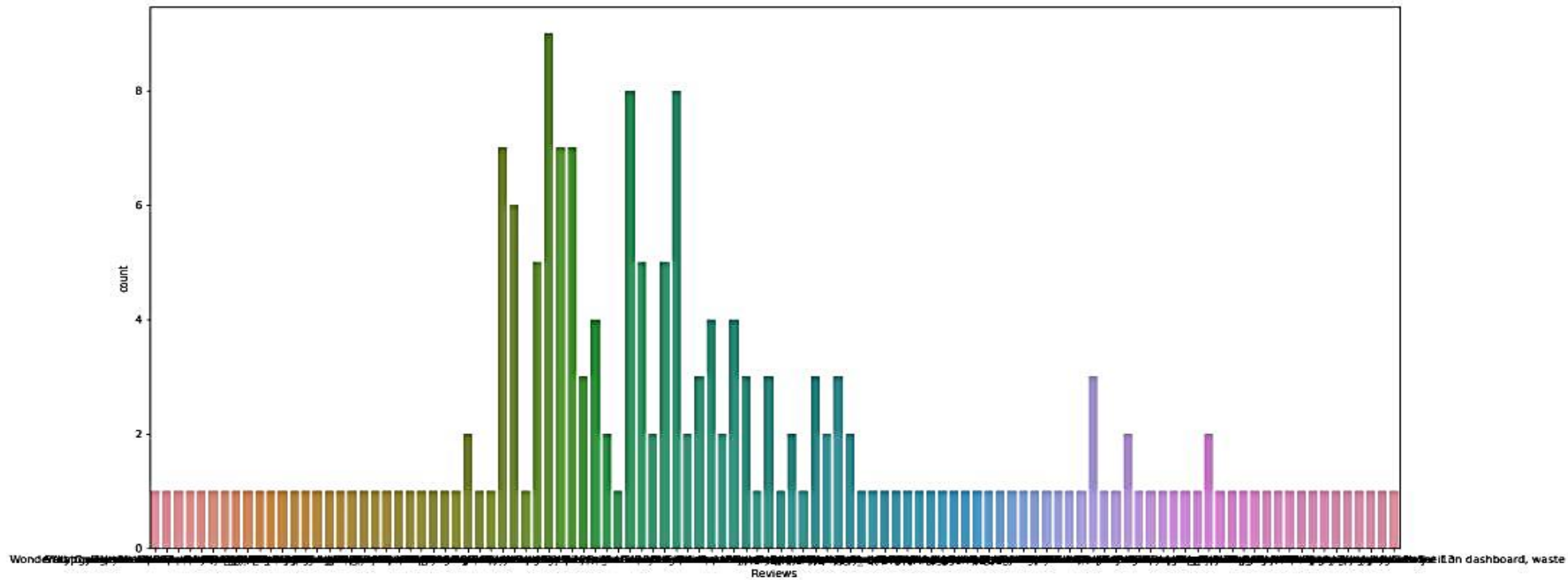
Positive review is on high as compared to others



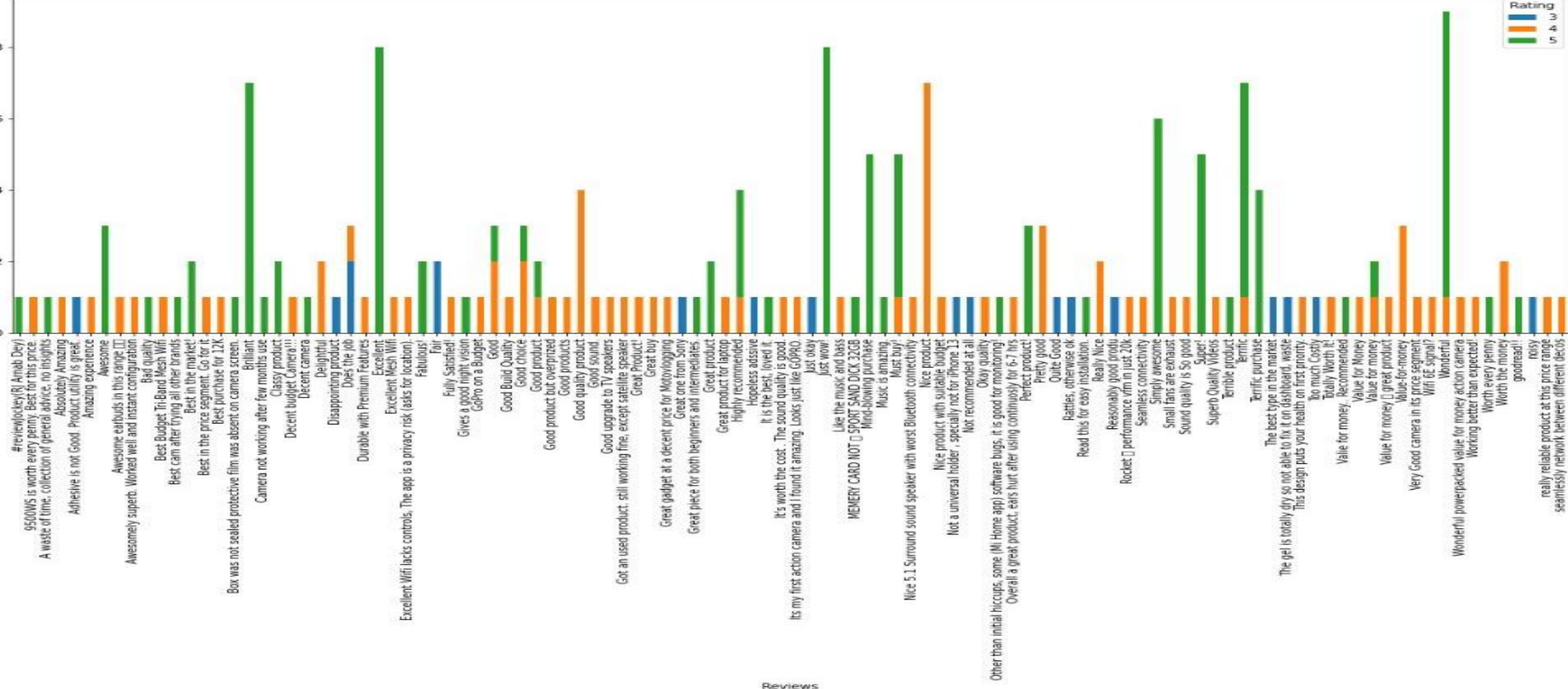
3. The count plot of the ratings of product:



4. The count plot of the reviews of product:



5.



- Data Sources and their formats

The data sources and their formats are from .csv file.

```
1 df=pd.read_csv('prediction excel sheet.csv')
2 df.head()
```

	Unnamed: 0	Rating	Reviews
0	0	4	Very Good camera in its price segment
1	1	4	Wonderful powerpacked value for money action c...
2	2	4	Great gadget at a decent price for Motovlogging
3	3	4	GoPro on a Budget
4	4	4	really reliable product at this price range

- Data Preprocessing Done

The steps followed for the cleaning of the data is Label Encoder after then importing preprocessing there transform the target columns into features then lastly it has to set/fir for the data frame.

Data Preprocessing

```
1 from sklearn import preprocessing
2 features=df.drop(['Unnamed: 0','Reviews'],axis=1)
3 target=df['Rating']
4 col_names=list(features.columns)
5 scaler=preprocessing.StandardScaler()
6 features=scaler.fit_transform(features)
7 features=pd.DataFrame(features,columns=col_names)
8 features.describe().T
```

	count	mean	std	min	25%	50%	75%	max
Rating	198.0	-3.364312e-16	1.002535	-2.170608	-0.635678	0.899252	0.899252	0.899252

- Data Inputs- Logic- Output Relationships

The relationships between inputs and outputs can be studied extracting weights of the trained model. Regression is that relationships between them can be blocky or highly structured based on the training data. It requires the data scientist to train the algorithm with both labeled inputs and desired outputs.

- State the set of assumptions (if any) related to the problem under consideration

Presumptions are by using regression label encoding, classifier, selection of the best models, confusion matrix that it means the relationship between the dependent and independent variables look fairly linear. Thus, our linearity assumption is satisfied.

- Hardware and Software Requirements and Tools Used
By importing many libraries are

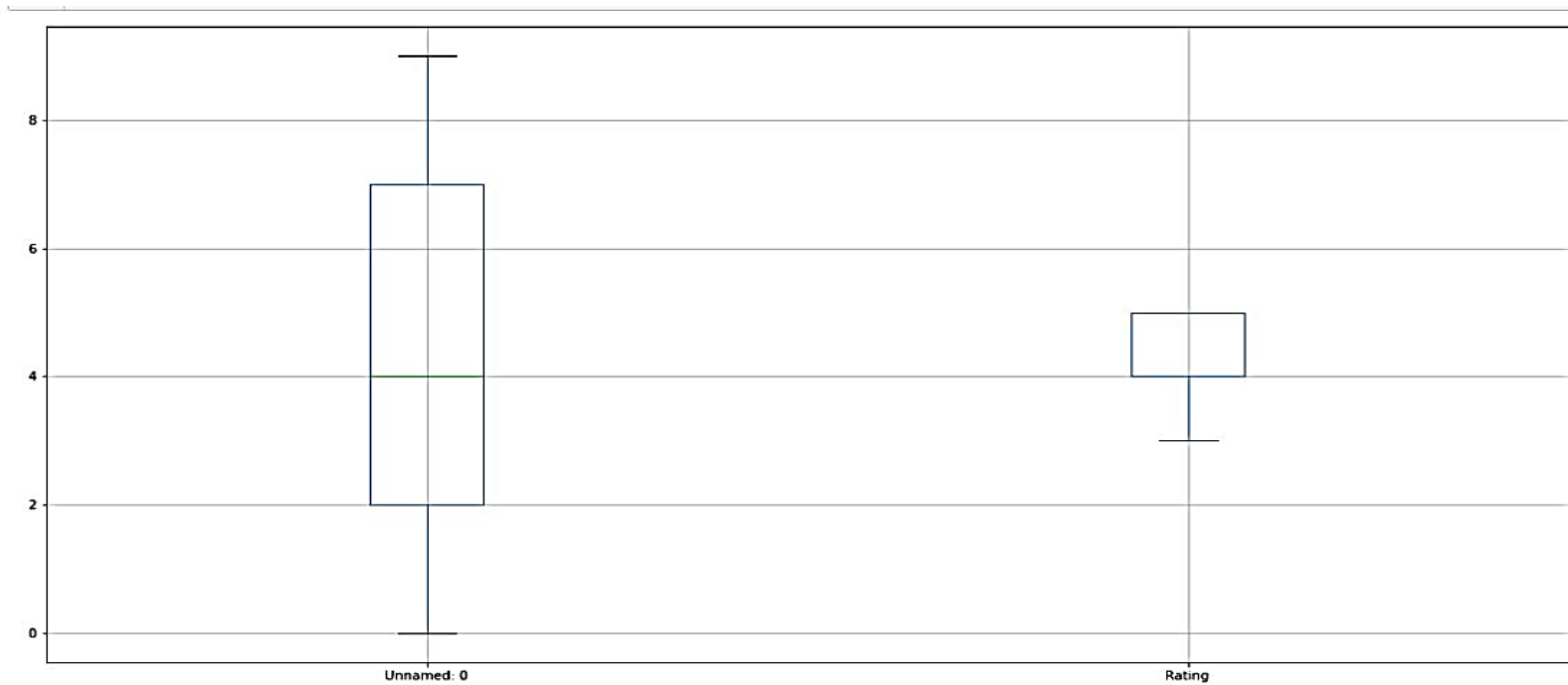
```
1 import pandas as pd
2 import seaborn as sns
3 import matplotlib.pyplot as plt
4 %matplotlib inline
5 import numpy as np
6 from sklearn.linear_model import LinearRegression
7 from sklearn.model_selection import train_test_split
8 import pickle
9 from sklearn.datasets import make_classification
10 from sklearn.linear_model import LogisticRegression
11 from sklearn.metrics import f1_score
12 from matplotlib import pyplot
13
14 import warnings
15 warnings.filterwarnings('ignore')
```

MODEL/s DEVELOPMENT AND EVALUATION

- Identification of possible problem-solving approaches(methods)
The collection and interpretation of data in order to uncover patterns and trends. It is a component of data analytics. Statistical analysis can be used in situations like gathering research interpretations, statistical modelling or designing surveys and studies. The approaches/methods of identification are descriptive and inferential statistics which are describes as the properties of sample and population data, and inferential statistics which uses those properties to test hypotheses and draw efficient conclusions in terms of outputs.

- Testing of Identified Approaches(Algorithms)

There is no outliers.



- Run and Evaluate selected models

DATA SCALING

```
1 from sklearn.preprocessing import StandardScaler
2 # Data Scaling Formula  $Z = (x - \text{mean}) / \text{std}$ 
3 scaler = StandardScaler()
4 X_scaled = scaler.fit_transform(X)
5 X_scaled
```

```
array([[ -1.52399933,  1.39744572],
       [ -1.17024567,  1.49836764],
       [ -0.81649202, -0.35186749],
       [ -0.46273837, -0.78919579],
       [ -0.10898471,  1.70021147],
       [  0.24476894,  0.28730465],
       [  0.59852226,  0.92647678],
       [  0.95227625, -1.22652409],
       [  1.3060299 ,  1.33016444],
       [  1.65978356, -0.04910174],
       [ -1.52399933, -1.52928984],
       [ -1.17024567,  1.73385211],
       [ -0.81649202, -0.95739898],
       [ -0.46273837, -1.59657112],
       [ -0.10898471,  1.43108636],
       [  0.24476894,  0.72463295],
       [  0.59852226, -1.69749303],
       [  0.95227625, -0.99103962],
       [  1.3060299 ,  0.75927262],
       [  1.65978356,  0.25927262]])
```

GRADIENT BOOSTING CLASSIFIER:

Accuracy score: Train result: 94.2%

Test result: 63.3%

```
1 # hyperparameter tuning
2 from sklearn.model_selection import GridSearchCV
3 grid_param = {
4     'max_depth': range(4,8),
5     'min_sample_split': range(2,8,2),
6     'learning_rate': np.arange(0.1,0.3)
7 }
```

=====Train Result=====

Accuracy Score: 94.20289855072464

CLASSIFICATION REPORT :

	3	4	5	accuracy	macro avg	weighted avg
precision	0.937500	0.962963	0.926471	0.942029	0.942311	0.943007
recall	1.000000	0.896552	0.969231	0.942029	0.955261	0.942029
f1-score	0.967742	0.928571	0.947368	0.942029	0.947894	0.941683
support	15.000000	58.000000	65.000000	0.942029	138.000000	138.000000

Confusion Matrix:

```
[[15 0 0]
 [1 52 5]
 [0 2 63]]
```

=====Test Result=====

Accuracy: 63.33333333333333

CLASSIFICATION REPORT :

	3	4	5	accuracy	macro avg	weighted avg
precision	0.0	0.526316	0.777778	0.633333	0.434698	0.646686
recall	0.0	0.454545	0.800000	0.633333	0.418182	0.633333
f1-score	0.0	0.487805	0.788732	0.633333	0.425512	0.638956
support	3.0	22.000000	35.000000	0.633333	60.000000	60.000000

Confusion Matrix:

```
[[0 2 1]
 [5 10 7]
 [0 7 28]]
```

Selecting the best model

```
: 1 print('The accuracy of the KNN Model is 0.6166666666666667')
   2 print('The accuracy of the Random Forest Model is 0.65')
```

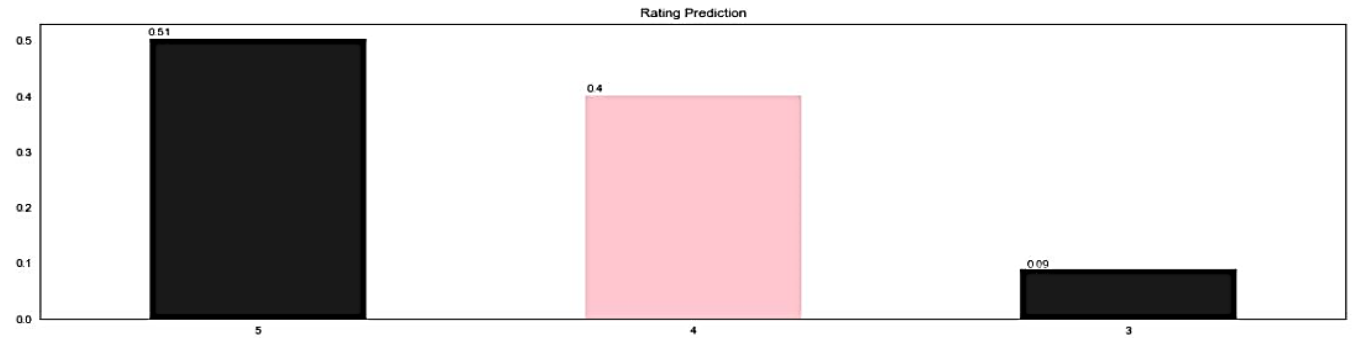
The accuracy of the KNN Model is 0.6166666666666667

The accuracy of the Random Forest Model is 0.65

- Visualizations

RESAMPLING THE LABEL

```
1 plt.figure(figsize=(20,5))
2 ax=df.Rating.value_counts(normalize=True).plot(kind='bar', color=['black', 'pink'], alpha=0.9, rot=0)
3 plt.title('Rating Prediction')
4 for i in ax.patches:
5     ax.annotate(str(round(i.get_height(),2)),(i.get_x() * 1.01, i.get_height() * 1.01))
6
7 plt.show()
```



USING THE HEATMAP FOR THE DATASET

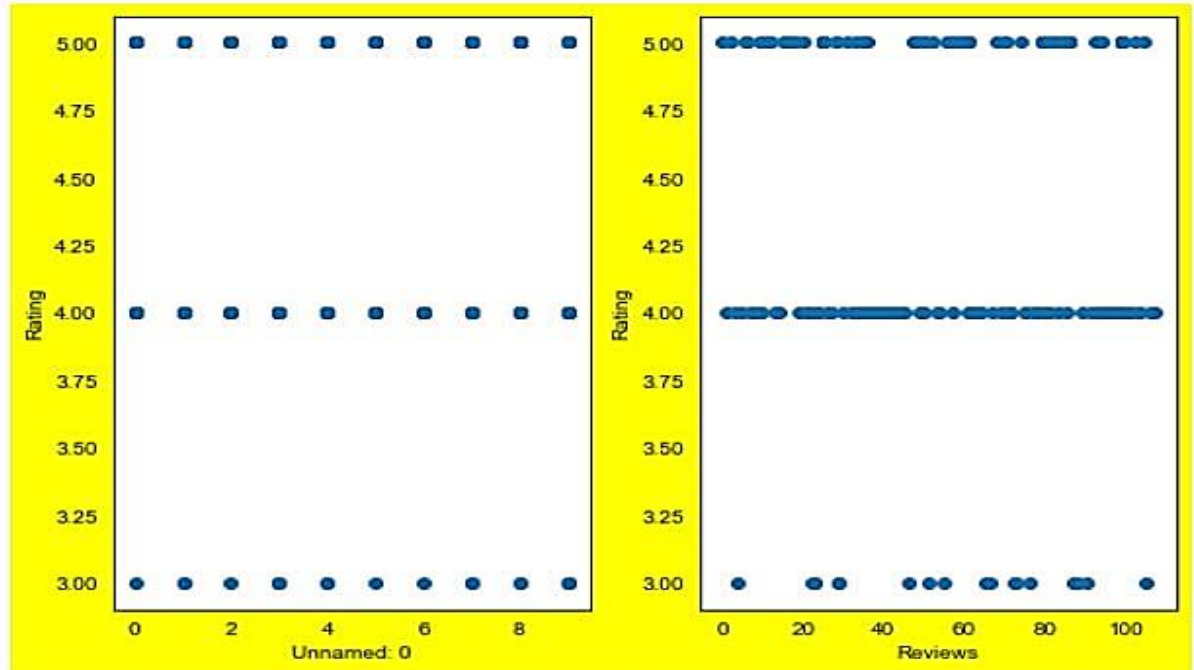
```
1 plt.figure(figsize=(20,5))
2 sns.heatmap(df.isnull(),cbar=False,cmap='twilight')
```

<AxesSubplot:>



VISUALIZING THE RELATIONSHIPS

```
1 # Visualizing relationship
2 plt.figure(figsize=(15,10), facecolor='yellow')
3 plotnumber = 1
4
5 for column in X:
6     if plotnumber<=8 :
7         ax = plt.subplot(2,4,plotnumber)
8         plt.scatter(X[column],y)
9         plt.xlabel(column,fontsize=10)
10        plt.ylabel('Rating',fontsize=10)
11        plotnumber+=1
12 plt.tight_layout()
```



The relationship between the dependent and independent variables look good in linear. Thus, our linearity assumption is satisfied.

- Interpretation of the Results

The results were interpreted from the visualizations, preprocessing and modelling:

1. The purpose of this case is to understand the rise of e-commerce which brought a significant rise in the importance of customer reviews and evaluate used to predict rating and to develop a strategy that utilizes data mining techniques towards ratings prediction.
2. There are two main methods to approach this problem. The first one is based on review text content analysis and uses the principles of natural language process (the NLP method). This method lacks the insights that can be drawn from the relationship between costumers and items.
3. The model of the independent variables and dependent variables are exactly vary with the variables.
4. It can accordingly manipulating the strategy of the areas that will yield high returns as it make easier for the clients.
5. By visualizations there are many things to be noted when it will according to work each other
6. By preprocessing the data it means that from the help of label encoder helps the dataset column to transform to fit another column in to it.

CONCLUSION

In conclusion, combining the formerly know data about each user' similarity to other users with the sentiment analysis of the rating and reviews itself, does help to improve the model prediction of rate the user's review, will get the purpose.

THANK YOU