# RATING PREDICTION PROJECT

# PROBLEM STATEMENT:

We have a client who has a website where people write different reviews for technical products. Now they are adding a new feature to their website i.e. The reviewer will have to add stars(rating) as well with the review. The rating is out 5 stars and it only has 5 options available 1 star, 2 stars, 3 stars, 4 stars, 5 stars. Now they want 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.

# UNDERSTANDING:

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.

# EDA STEPS AND VISUALIZATIONS:

## 1. Data Collection

```
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 |

# 2. Data Cleaning

```
1  total=df.isnull().sum().sort_values(ascending=False)
2  percent = (df.isnull().sum()/df.isnull().count()).sort_values(ascending=False)
3  missing = pd.concat([total, percent], axis=1, keys=['Total' , 'Percent'])
4  missing.head()
```

|  | Total | Percent |
|---|---|---|
| Unnamed: 0 | 0 | 0.0 |
| Rating | 0 | 0.0 |
| Reviews | 0 | 0.0 |

# 3. Univariate Analysis

```python
from sklearn.preprocessing import StandardScaler
# Data Scaling Formula Z=(x-mean)/std
scaler = StandardScaler()
X_scaled=scaler.fit_transform(X)
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],
       [ 0.5985226 ,  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],
       [ 0.5985226 , -1.69749303],
       [ 0.95227625, -0.99103962],
```
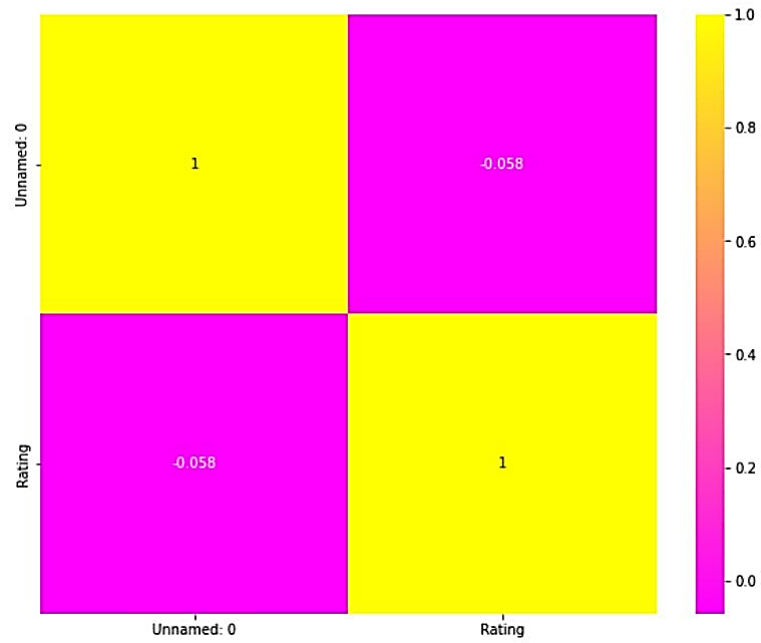
Rating

Reviews

```
1  plt.figure(figsize=(10,8))
2  sns.heatmap(corr,annot=True,cmap='spring')
3  plt.show()
```
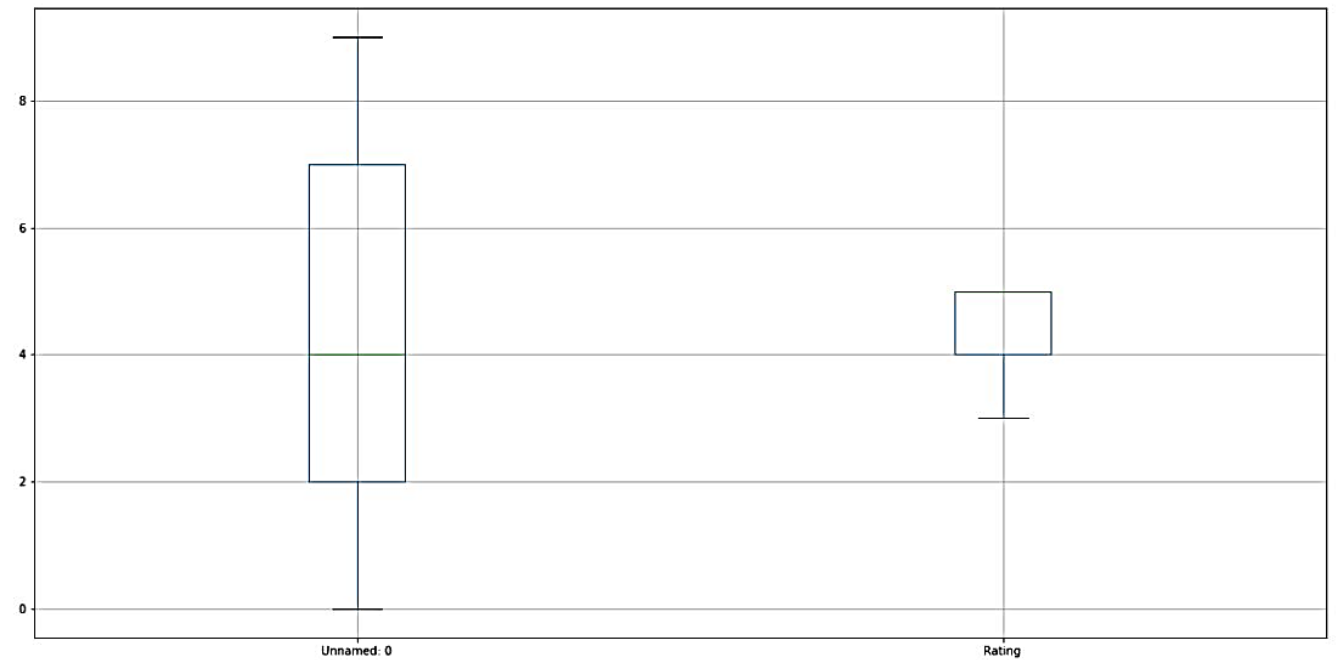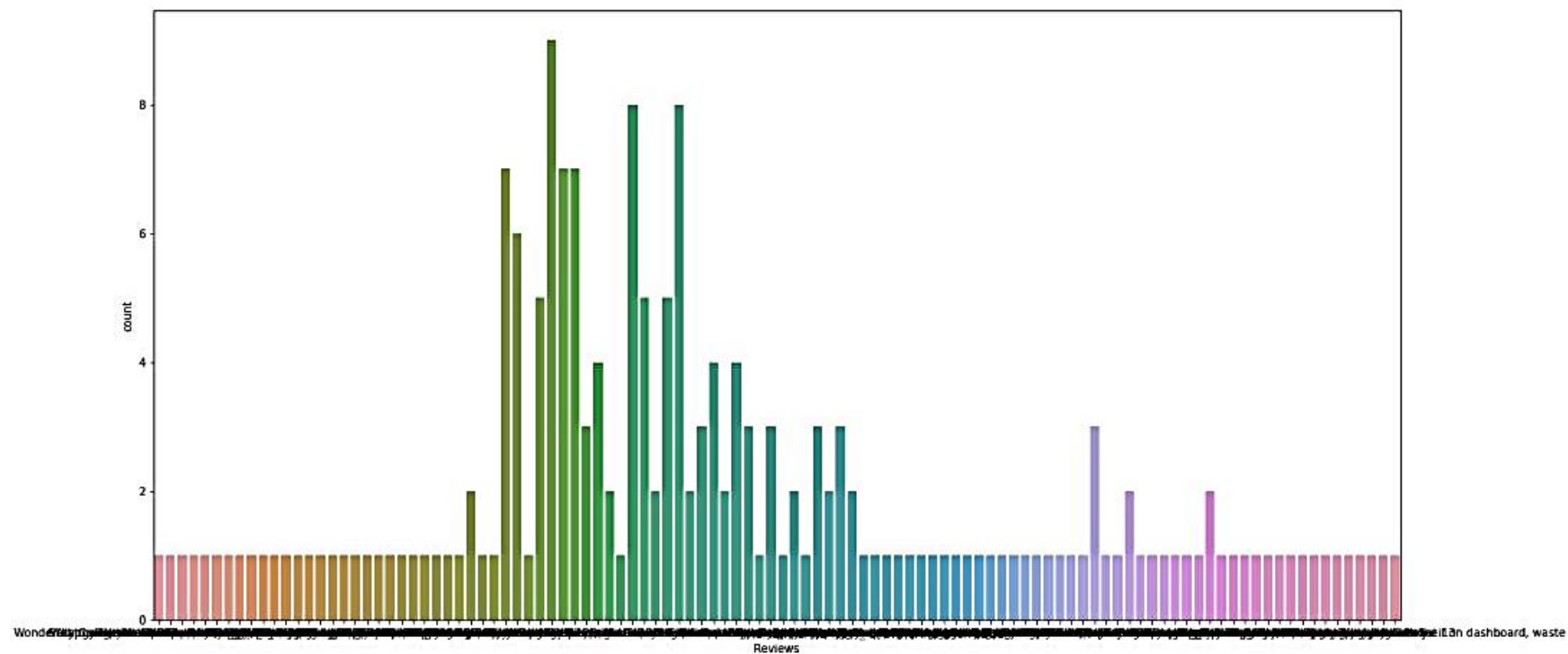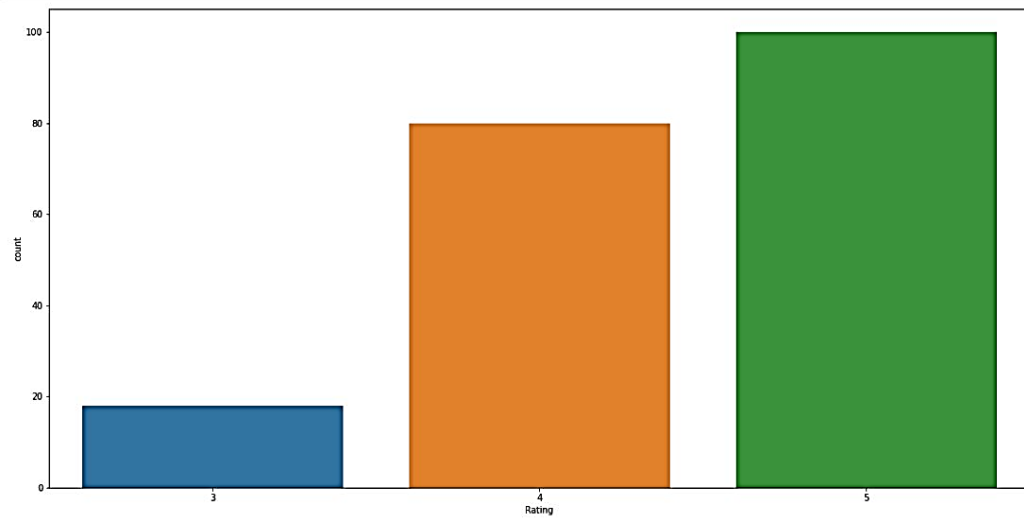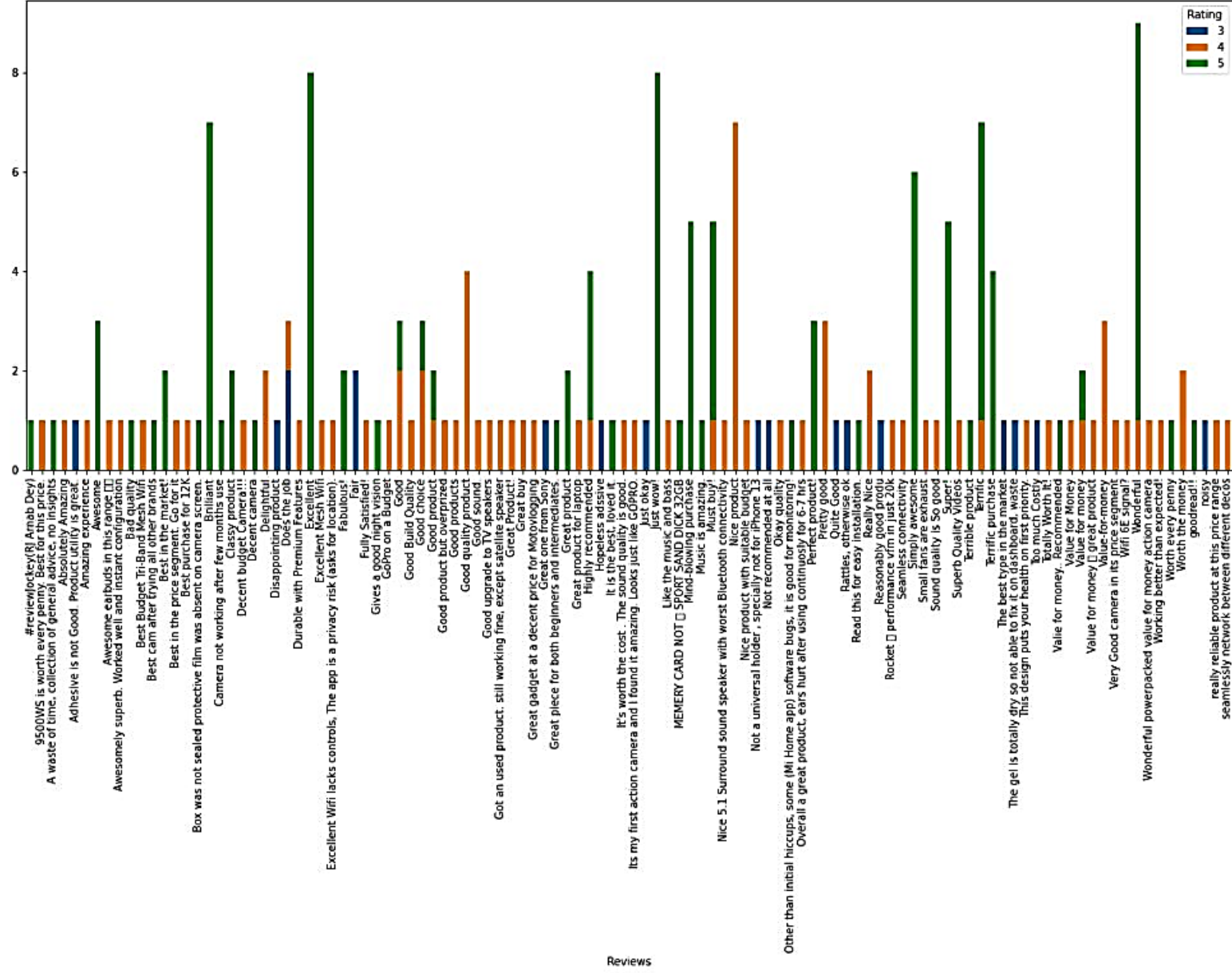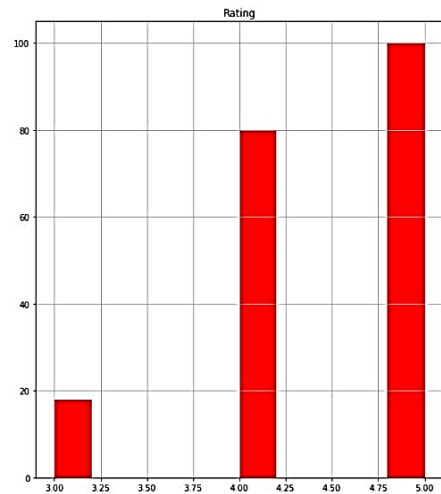
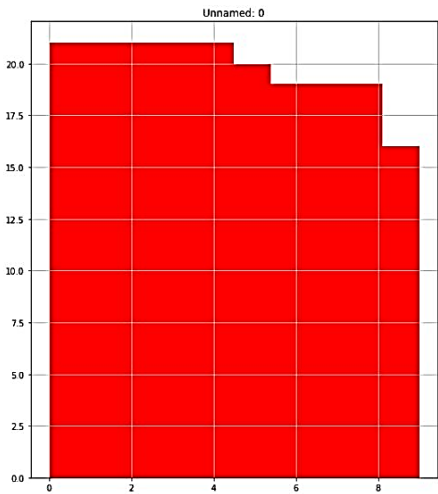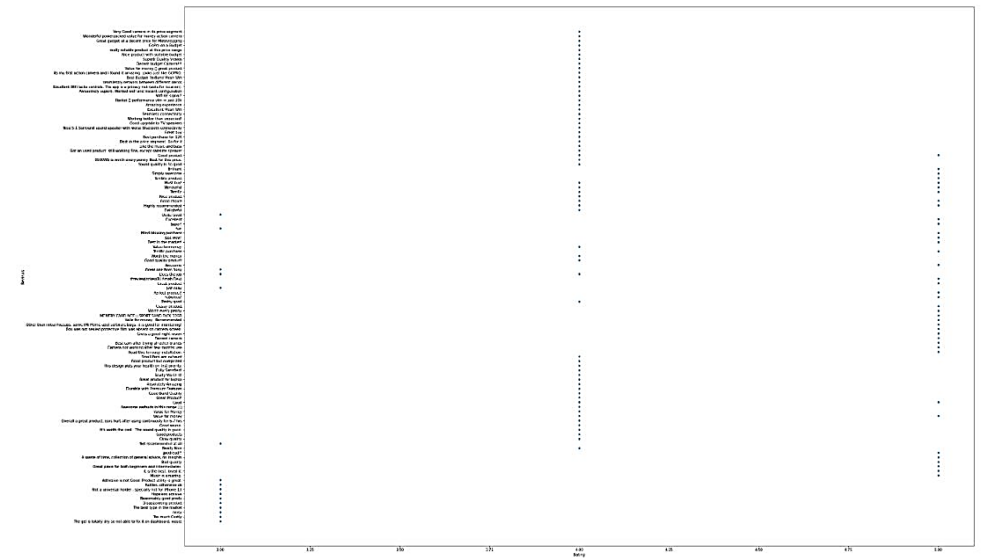

```
1  plt.figure(figsize=(20,10))
2  df.boxplot();
```

```
1  df.hist(color='r',figsize=(20,10))
2  plt.show()
```
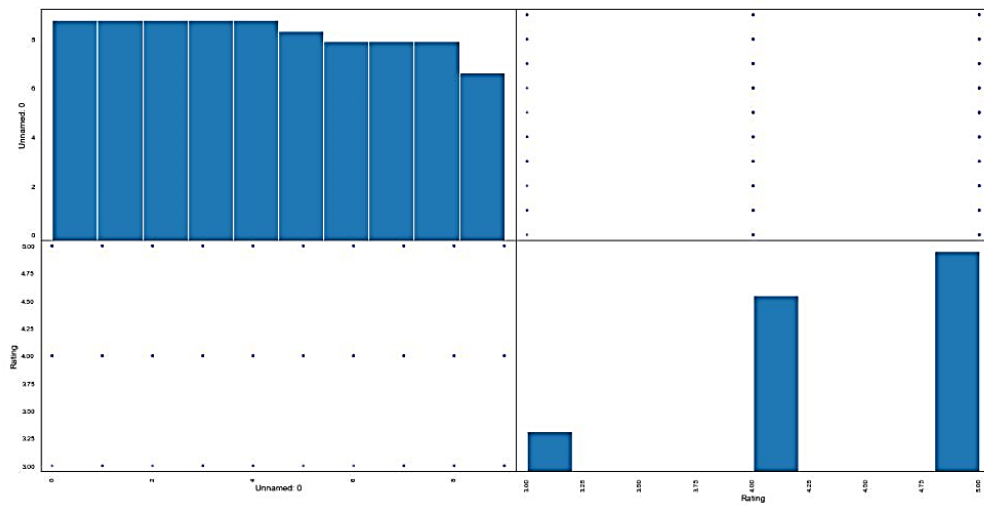


```
1  plt.figure(figsize=(40,28))
2  sns.scatterplot(df['Rating'],df['Reviews'])
3  print(df[['Rating','Reviews']].corr())
4  plt.show()
```

```
          Rating
Rating    1.0
```
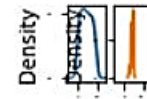


```
1  from pandas.plotting import scatter_matrix
2  scatter_matrix(df,figsize=(20,10),color='b')
3  plt.show()
```



```
1  df.plot(kind='density',subplots=True,layout=(5,15),legend=False,fontsize=2);
```
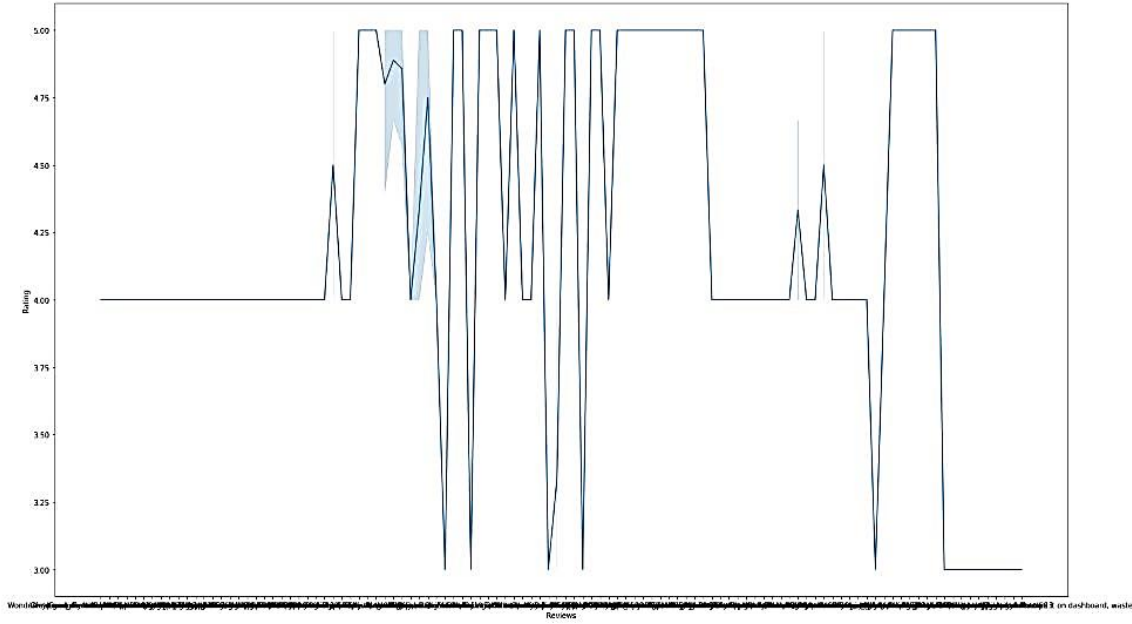
```
1  plt.figure(figsize=(25,15))
2  sns.lineplot(x='Reviews',y='Rating',data=df)
```
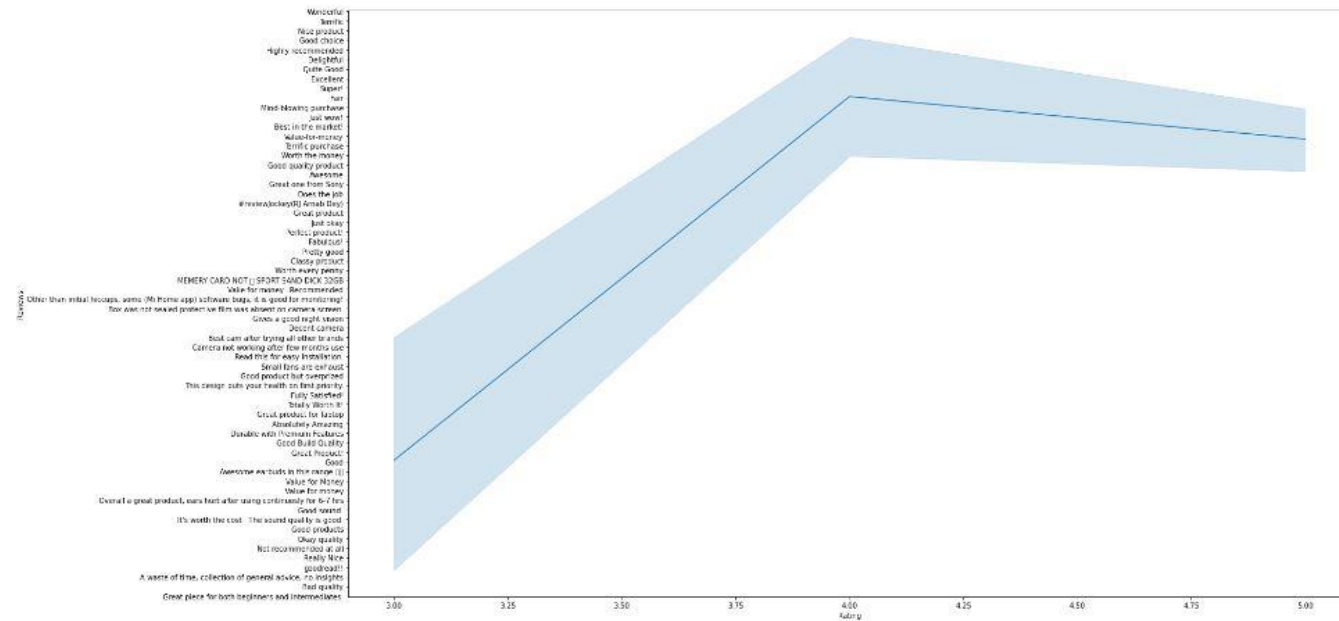
<AxesSubplot:xlabel='Reviews', ylabel='Rating'>



```
1  plt.figure(figsize=(25,15))
2  sns.lineplot(x='Rating',y='Reviews',data=df)
```

<AxesSubplot:xlabel='Rating', ylabel='Reviews'>

```
1  plt.figure(figsize=(40,28))
2  sns.relplot(x='Rating', y='Reviews',data=df)
3  plt.xticks(rotation=90);
```

<Figure size 2880x2016 with 0 Axes>



```
1  def correlation_heatmap(df):
2      ax=plt.subplots(figsize=(20,10))
3      colormap=sns.diverging_palette(220, 10, as_cmap = True)
4      ax=sns.heatmap(df.corr(),cmap="magma",annot=True,linewidths=0.1,vmax=1.0,linecolor='white',annot_kws={'fontsize':12})
5      plt.title('Correlation',y=1.05,size=15)
6
7  correlation_heatmap(df)
```

Rating Prediction

```
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()
```
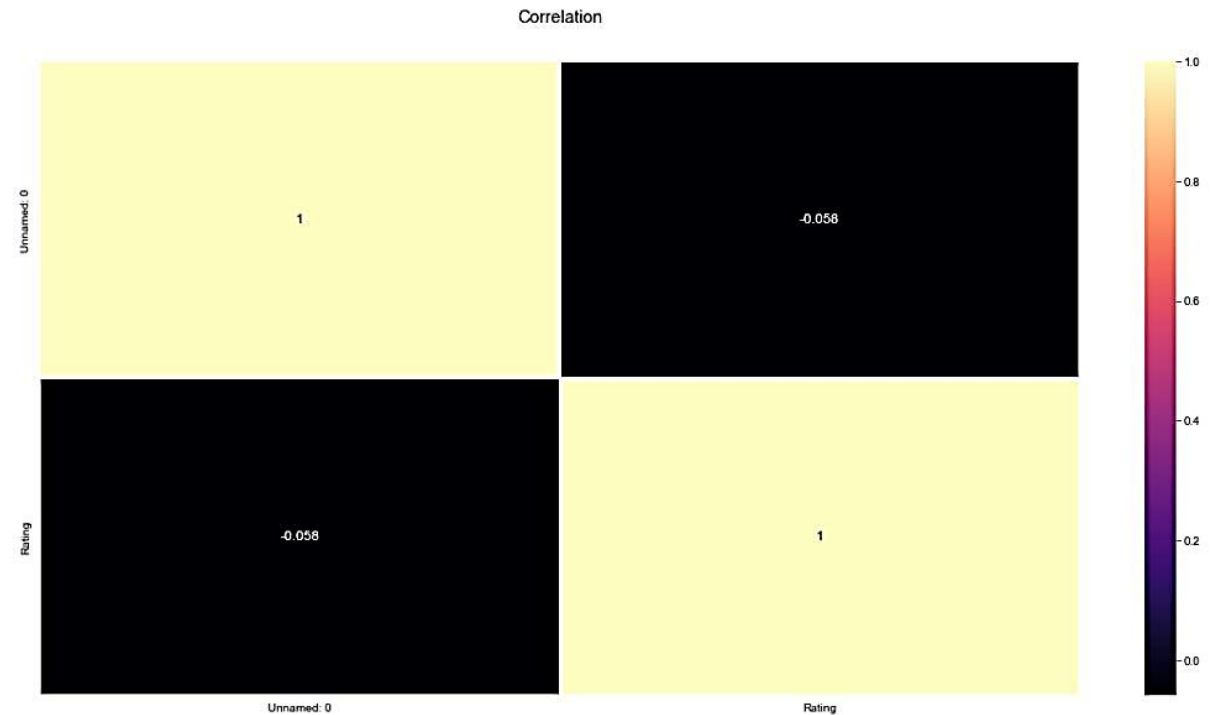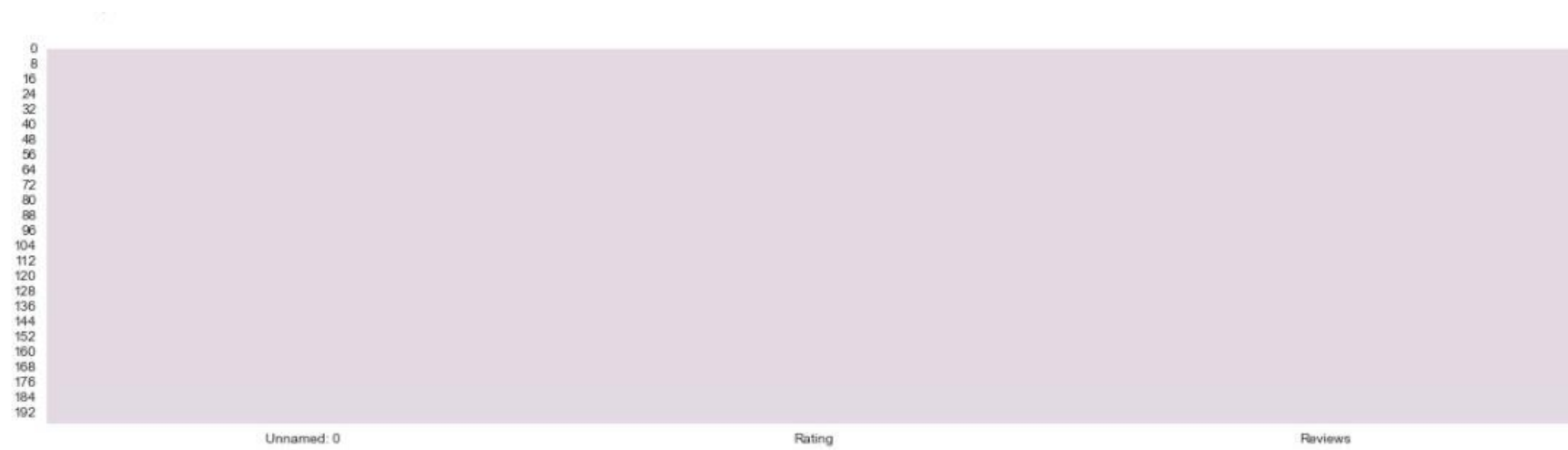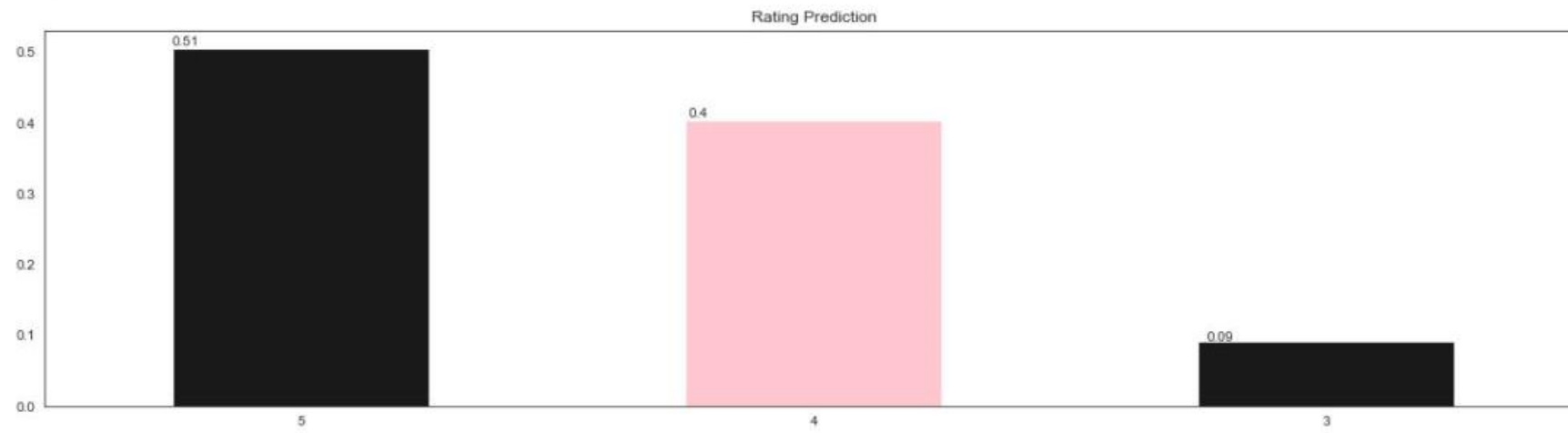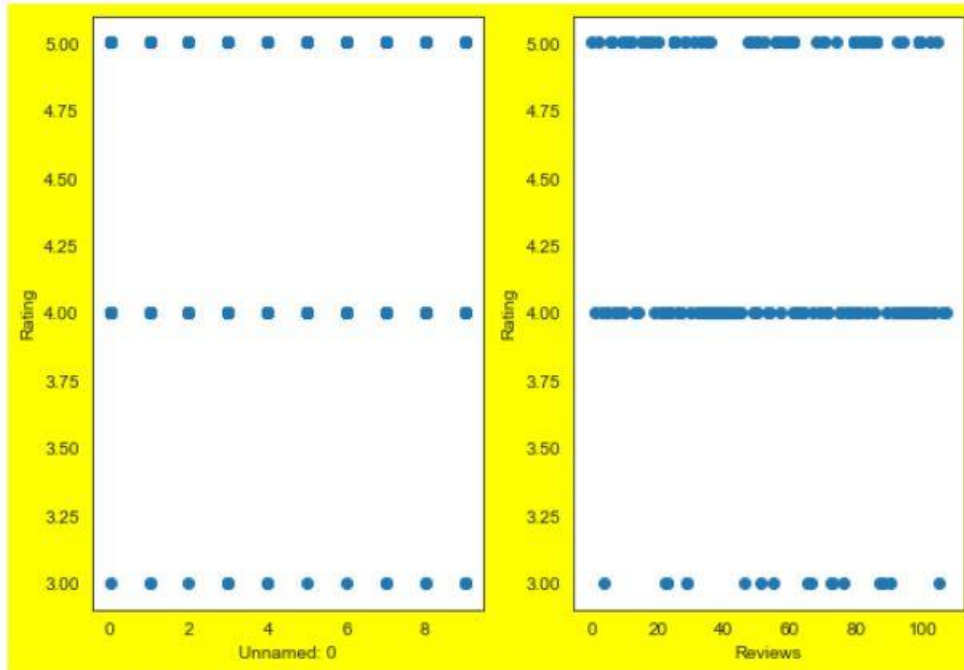


**VISUALIZING THE RELATIONSHIP WITH RATING**

```
1  from sklearn.preprocessing import StandardScaler
2  # Data Scaling Formula Z=(x-mean)/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],
       [ 0.5985226 ,  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],
       [ 0.5985226 , -1.69749303],
       [ 0.95227625, -0.9910396.]
```

```
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  }
```

# STEPS AND ASSUMPTIONS TO COMPLETE THE PROJECT:

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.

# FINALIZED MODEL:

Accuracy score:
Train Result : 94.2%
Test Result : 63.3%



```
=====================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

Cofusion 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]]
```

# 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