

RATINGS PREDICTION PROJECT

Submitted by: NIPAM GOGOI

ACKNOWLEDGMENT I would like to thank Flip Robo Technologies for providing me with the opportunity to work on this project from which I have learned a lot. I am also grateful to Miss Khushboo Garg for her constant guidance and support.

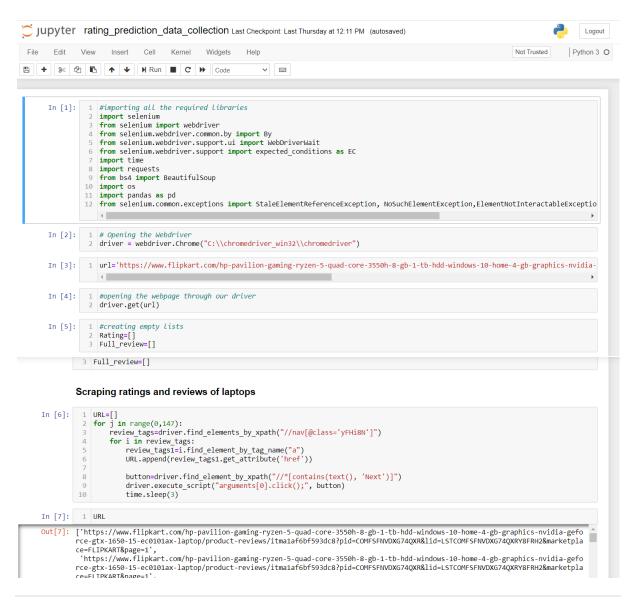
INTRODUCTION

BUSINESS PROBLEM FRAMING

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 rating. So we, we have to build an application which can predict the rating by seeing the review.

ANALYTICAL PROBLEM FRAMING

DATA COLLECTION - WEBSCRAPPING



```
Out[16]: 1410
      In [17]: 1 len(Full_review)
 Out[17]: 1410
      In [18]: 1 url='https://www.flipkart.com/apple-macbook-air-core-i5-5th-gen-8-gb-128-gb-ssd-mac-os-sierra-mqd32hn-a/product-reviews/itm0
                   #opening the webpage through our driver
driver.get(url)
      In [19]:
                  for i in range(0,328):
for j in driver.find_elements_by_xpath("//div[@class='_3LWZlK _1BLPMq' or @class='_3LWZlK _32lA32 _1BLPMq' or @class='_3
      In [20]:
                                Ratings.append(j.text)
                           for k in driver.find_elements_by_xpath("//div[@class='t-ZTKy']/div/div"):
    Full_review.append(k.text.replace('\n',' '))
                           time.sleep(1)
button=driver.find_element_by_xpath("//*[contains(text(), 'Next')]")
driver.execute_script("arguments[0].click();", button)
                           time.sleep(2)
                  13
      In [21]: 1 len(Ratings)
Out[124]:
                     0
                                            Its an absolute beast if u know what are the n...
                     1
                                               This is the best laptop in this range.I reciev...
                     2
                               5
                                            Good product as used of now.... Everything is ...
                                5 AWESOME LAPTOP. It supports many high spec gam...
                                               For that price... it's exceptionally good. Pla...
                39364
                                5
                                        This device completely covers my bunglow.. Goo...
                39365
                                               a great product from tenda,....i literally kee..
                39366
                                3
                                               Good router bt Flipkart fluctuates it's price....
                                4
                                             Too much gd...lokking nice and fit..range is a...
                                          Really great product, delivery was awesome. Al...
                39368
               39369 rows × 2 columns
                 #converting into csv format
Reviews.to_csv("Rating_Prediction_data.csv")
In [125]:
In [126]:
                 1 #closing the driver
```

2 driver.close()

MODEL BUILDING AND EVALUATION

```
Jupyter Ratings_prediction_project_Model_Building Last Checkpoint: 7 hours ago (autosaved)
                                                                                                                                                     Logo
             View Insert Cell Kernel Widgets Help
                                                                                                                                                  Python 3
# importing all necessery modules
import pandas as pd
from sklearn.feature_extraction.text import CountVectorizer,TfidfVectorizer,TfidfTransformer
from sklearn.ensemble import RandomForestClassifier ,AdaBoostClassifier
from sklearn.linear_model import togisticRegression
from sklearn.metrics import precision_score, recall_score,accuracy_score, classification_report, f1_score ,confusion_matrix
from sklearn.model_selection import train_test_split
from sklearn.model_selection import train_test_split
      In [1]:
                    from wordcloud import WordCloud, STOPWORDS
                  9 import matplotlib.pyplot as plt
      In [2]: 1 df = pd.read_csv("Rating_Prediction_data.csv")
      In [5]: 1 df.columns
      Out[5]: Index(['Unnamed: 0', 'Ratings', 'Full_review'], dtype='object')
      In [6]: 1 df = df.drop("Unnamed: 0", axis=1)
          2 def Word_Cloud(str_List):
                   comment_words =
          4
                   stopwords = set(STOPWORDS)
                   # iterate through the csv file
          6
                   for val in str_List:
          8
         9
                         # typecaste each val to string
                         val = str(val)
         10
         11
                         # split the value
         12
         13
                         tokens = val.split()
         14
         15
                         # Converts each token into lowercase
                         for i in range(len(tokens)):
         16
         17
                               tokens[i] = tokens[i].lower()
         18
                         comment_words += " ".join(tokens)+" "
         19
         20
         21
         22
                   wordcloud = WordCloud(width = 800, height = 800,
                                          background_color ='white',
stopwords = stopwords,
         23
         24
         25
                                          min_font_size = 10).generate(comment_words)
         26
         27
                    # plot the WordCloud image
                   plt.figure(figsize = (8, 8), facecolor = None)
         28
                   plt.imshow(wordcloud)
         29
                   plt.axis("off")
         30
                   plt.tight_layout(pad = 0)
         31
         32
         33
                   plt.show()
```

```
In [45]: 1 df.Ratings.unique()
Out[45]: array([5, 4, 3, 1, 2], dtype=int64)
print()
       RATINGS : 5
       range
        fast working love find purchase
 3 y = df.Ratings
 In [9]: 1 # Splitting the train_df into training and test dataset
2 X_train,X_test,y_train,y_test=train_test_split(X,y,test_size=0.2,random_state=88)
x_test = vectorizer.transform(X_test)
In [80]: 1 # Random Forest Model
In [11]: 1 rand = RandomForestClassifier(n_estimators=100,criterion='entropy',max_features=None,class_weight='balanced')
         2 rand.fit(x_train, y_train)
Out[11]: RandomForestClassifier(class_weight='balanced', criterion='entropy',
                         max_features=None)
```

In [12]: 1 prediction = rand.predict(x_test)

```
CONFUSION MATRIX
        [[ 732 124
                   53
                       64 111]
          154
               47
                  57
                      41
                           51]
        [ 111
               61 111 178 190]
          50
               37
                  113
                      580
                          638]
               64 146 1230 2809]]
        ACCURACY
        0.5434340868681737
        REPORT
                    precision
                              recall f1-score
                                              support
                1
                       0.63
                               0.68
                                       0.65
                                               1084
                2
                       0.14
                               0.13
                                       0.14
                                                350
                3
                       0.23
                               0.17
                                       0.20
                                                651
                                       0.33
                                               1418
                4
                       0.28
                               0.41
                       0.74
                               0.64
                                               4371
                                       0.69
                                       0.54
                                                7874
           accuracy
                       9.49
                               0.41
                                                7874
          macro avg
                                       9.49
       weighted avg
                       0.57
                               0.54
                                       0.55
                                               7874
```

The Accuracy is very less. So Transforming into a binary classification Problem with only High Rating and Low Rating

```
In [14]:
            1 Target = list(df.Ratings)
In [63]:
               out = []
            1
               for i in Target:
            3
                    if i == 5 or i == 4:
                        out.append(1)
                    elif i == 1 or i == 2:
            5
            6
                        out.append(0)
            7
                    else:
                        out.append("Neutral")
            8
           1 df["Target"] = out
In [64]:
           1 df.Target.unique()
In [79]:
Out[79]: array([1, 'Neutral', 0], dtype=object)
In [79]: 1 df.Target.unique()
Out[79]: array([1, 'Neutral', 0], dtype=object)
       for i in list(df.Target.unique()):
    if i == 1:
In [76]:
                print("High Rated Reviews")
f i == 0:
                print("Worst Rated Reviews")
             print("Neutral Reviews")
Word_cloud(list(df[df.Target == i].sample(n=1500).Full_review))
         thing
                                                   buy
           great
                                         excellentphone
                                                            dslr
                              one
              display
                                                           get
         superb
```

```
In [65]: 1 df2 = df[df.Target != "Neutral"]
            1 # Defining Input data and Output to be predicted
2 X2 = df2.Full_review
3 y2 = list(df2.Target)
 In [66]:
In [67]: 1 # Splitting the train_df into training and test dataset
2 X_train2,X_test2,y_train2,y_test2=train_test_split(X2,y2,test_size=0.2,random_state=88)
4 x_test2 = vectorizer.transform(X_test2)
 In [69]: 1 rand = RandomForestClassifier(n_estimators=100,criterion='entropy',max_features=None,class_weight='balanced')
             2 rand.fit(x_train2, y_train2)
 Out[69]: RandomForestClassifier(class_weight='balanced', criterion='entropy',
                                    max_features=None)
 In [70]: 1 prediction = rand.predict(x_test2)
In [71]: 1 print('\n','CONFUSION MATRIX','\n',confusion_matrix(y_test2, prediction))
2 print('\n','ACCURACY','\n',accuracy_score(y_test2, prediction))
3 print('\n','REPORT','\n',classification_report(y_test2,prediction))
            CONFUSION MATRIX
            [[1127 324]
              248 5538]]
In [70]:
               1 prediction = rand.predict(x_test2)
                print('\n','CONFUSION MATRIX','\n',confusion_matrix(y_test2, prediction))
print('\n','ACCURACY','\n',accuracy_score(y_test2, prediction))
print('\n','REPORT','\n',classification_report(y_test2,prediction))
In [71]:
                CONFUSION MATRIX
                [[1127 324]
                [ 248 5538]]
                ACCURACY
                0.9209617244714661
                REPORT
                                      precision
                                                         recall f1-score
                                                                                       support
                               a
                                           0.82
                                                           0.78
                                                                           0.80
                                                                                          1/151
                               1
                                            0.94
                                                           0.96
                                                                           0.95
                                                                                          5786
                                                                           0.92
                                                                                          7237
                    accuracy
                   macro avg
                                           0.88
                                                           0.87
                                                                           0.87
                                                                                          7237
              weighted avg
                                           0.92
                                                           0.92
                                                                           0.92
                                                                                          7237
```

```
1 TORTSCICKER...ITC(X_CL.aTLLS, A_CL.aTLLS)
Out[72]: LogisticRegression(class weight='balanced', max iter=1000, random_state=5,
                                   solver='liblinear', tol=0.001)
In [73]: 1 prediction2 = logisticRegr.predict(x_test2)
            print('\n','CONFUSION MATRIX','\n',confusion_matrix(y_test2, prediction2))
print('\n','ACCURACY','\n',accuracy_score(y_test2, prediction2))
print('\n','REPORT','\n',classification_report(y_test2,prediction2))
             CONFUSION MATRIX
             [[1263 188]
             [ 449 5337]]
             ACCURACY
             0.9119801022523145
                              precision
                                              recall f1-score
                                                                     support
                                   0.74
                                               0.87
                                                            0.80
                                                                        1451
                                               0.92
                                                            0.94
                                                                        5786
                                                            0.91
                                                                        7237
                accuracy
                                   0.85
                                               0.90
                                                            0.87
               macro avg
                                                                        7237
           weighted avg
                                   0.92
                                                0.91
                                                            0.91
                                                                        7237
 In [ ]: 1
```

CONCLUSION

- The Data was collected using webscrapping from Flipkart for various products, like Laptop, mobiles etc. There are 2 columns in the Dataset: Ratings and Reviews.
- EDA was done on the dataset and Word Clouds were created for all types of Reviews.
- The ML models were not performing well for predicting all the ratings. Since its
 hard to distinguish between neutral reviews and good or bad reviews so results
 in inaccurate prediction of Ratings.
- Hence transformed into binary classification problem, taking 1 or 2 ratings as Bad ratings and 4 or 5 as good ratings. The final models performed well for this set up with around 91-92 % accuracy.
- We can directly rate the cases to predicted as 5, 3 or 1 on the unseen data as ratings based on review.