



Ratings Prediction Project

Submitted by:
SONALI WARKHADE

ACKNOWLEDGMENT

I would like to express my thanks of gratitude to mentors Shubham Yadav as well as Flip Robo who gave me the golden opportunity to do this wonderful project on the topic Micro Credit Defaulter, which also helped me in doing a lot of research and I came to know about so many new things. I am really very thankful to them.

I am making this project to increase my knowledge.

INTRODUCTION

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.

This project consists on of 2 phases:

Data Collection Phase

We need to scrape the reviews of different laptops, Phones, Headphones, smart watches, Professional Cameras, Printers, Monitors, Home theater, Router from different e-commerce websites. Basically, we need these columns -

- 1) reviews of the product.
- 2) rating of the product.

Model Building Phase

After collecting the data, you need to build a machine learning model. Before model building do all data pre-processing steps involving NLP. Try different models with different hyper parameters and select the best model.

Analytical Problem Framing

With the help of Selenium I have scrapped data from different websites like, Snapdeal, Flipkart and Amazon and used Pandas library to save the data in excel file. Just taking a glance on basic code for scrapping from different websites.

```
1 # Opening the homepage of Nykaa
2 url = "https://www.nykaa.com/"
3 driver.get(url)
4
5 time.sleep(2)
6
7 search_bar = driver.find_element_by_class_name("css-1p0hsio") # Locating searc_bar by id
8 search_bar.clear() # clearing search_bar
9 search_bar.send_keys('Headphones') # sending user input to search bar
10 search_bar.submit()
```

```
1 links1=[]
2 driver.execute_script("window.scrollTo(0,document.body.scrollHeight)")
3 time.sleep(2)
4 try:
5     for l in driver.find_elements_by_xpath('//div[@class="product-list-box card desktop-cart"]/a'):
6         links1.append(l.get_attribute('href'))
7         time.sleep(1)
8 except NoSuchElementException:
9     links1.append("-")
10    time.sleep(1)
```

```
1 len(links1)
```

```
1 Rating1=[]
2 Review1=[]
3 Review_title1=[]
4 for url in links1:
5     driver.get(url)
6     time.sleep(2)
7
8     try:
9         all_reviews = driver.find_element_by_xpath('//a[@class="all-review-btn"]')
10        a =all_reviews.get_attribute('href')
11        driver.get(a)
12
13        dropdown = driver.find_element_by_xpath('//*[id="js-pd-scroll-start"]/div[1]/div[2]/section/div/div')
14        negative = driver.find_element_by_xpath('//*[id="js-pd-scroll-start"]/div[1]/div[2]/section/div/div/span')
15        action= ActionChains(driver)
16        action.move_to_element(dropdown).move_to_element(negative).click().perform()
17
18        # Extracting stars
19        try:
20            rating = driver.find_elements_by_xpath('//div[@class="desc-section"]/div[1]/div[1]')
21            for r in rating:
22                Rating1.append(r.text)
23        except NoSuchElementException:
24            Rating1.append('-')
25
26        time.sleep(1)
27        # Extracting review
28        try:
29            review = driver.find_elements_by_xpath('//div[@class="review-details"]/section/p')
30            for r in review:
31                Review1.append(r.text)
32        except NoSuchElementException:
33            Review1.append('-')
34
35        time.sleep(1)
36        # Extracting review
37        try:
38            review_title = driver.find_elements_by_xpath('//div[@class="review-details"]/section/h4')
39            for r in review_title:
40                Review_title1.append(r.text)
41        except NoSuchElementException:
42            Review_title1.append('-')
43    except TimeoutException:
44        pass
45    except NoSuchElementException:
46        pass
```

Saving the scrapped data into a csv file:

```
1 # dictionary of lists
2 dict = {'Rating': Rating1, 'Review': Review1, 'Review_title': Review_title1}
3
4 df = pd.DataFrame(dict)
5
6 # saving the dataframe
7 df.to_csv('Nykaa.csv')
```

Snapdeal

```
1 items=['Headphones', 'Phones', 'smart watches', 'Monitors', 'Home theater', 'Router']
2 links2=[]
3 for i in items:
4
5     # Opening the homepage of snapdeal
6     url = "https://www.snapdeal.com/"
7     driver.get(url)
8
9     time.sleep(2)
10
11     search_bar = driver.find_element_by_xpath('//*[@id="inputValEnter"]') # Finding the search bar using it's xpath
12     search_bar.clear() # clearing search_bar
13     search_bar.send_keys(i) # sending user input to search bar
14     search_button = driver.find_element_by_xpath('//*[@id="sdHeader"]/div[4]/div[2]/div/div[2]/button') # Locating sea
15     search_button.click()
16
17     try:
18         for l in driver.find_elements_by_xpath('//*[@class="dp-widget-link noUdLine hashAdded"]'):
19             links2.append(l.get_attribute('href'))
20             time.sleep(1)
21     except NoSuchElementException:
22         links2.append("-")
23         time.sleep(1)
```

```
1 len(links2)
```

240

```
1 Rating2=[]
2 Review2=[]
3 Review_title2=[]
4 for url in links2:
5     driver.get(url)
6     time.sleep(2)
7
8     # Extracting stars
9     try:
10         rating = driver.find_elements_by_xpath('//div[@class="user-review"]//div[@class="rating"]')
11         for r in rating:
12             child= r.find_elements_by_tag_name("i")
13             Rating2.append(len(child))
14     except NoSuchElementException:
15         Rating2.append("-")
16
17     time.sleep(1)
18     # Extracting review
19     try:
20         review = driver.find_elements_by_xpath('//div[@class="commentreview"]/div/div[2]/div[2]')
21         for r in review:
22             Review2.append(r.text)
23     except NoSuchElementException:
24         Review2.append("-")
25
26     time.sleep(1)
27     # Extracting review
28     try:
29         review_title = driver.find_elements_by_xpath('//div[@class="commentreview"]/div/div[2]/div[1]/div[2]')
30         for r in review_title:
31             Review_title2.append(r.text)
32     except NoSuchElementException:
33         Review_title2.append("-")
```

```
1 # dictionary of lists
2 dict2 = {'Rating': Rating2, 'Review': Review2, 'Review_title': Review_title2}
3
4 df = pd.DataFrame(dict2)
5
6 # saving the dataframe
7 df.to_csv('Snapdeal.csv')
```

```

1 # Opening the homepage of flipkart
2 url = "https://www.flipkart.com/"
3 driver.get(url)
4
5 try:
6     login_X_button = driver.find_element_by_xpath('//button[@class="_2KpZ6l _2doB4z"]') # Button to close login popup
7     login_X_button.click()
8 except NoSuchElementException :
9     print("No Login page")
10
11 time.sleep(2)
12 try:
13     search_bar = driver.find_element_by_xpath('//*[@id="container"]/div/div[1]/div[1]/div[2]/div[2]/form/div/div/input')
14     search_bar.clear() # clearing search_bar
15     search_bar.send_keys('laptops') # sending user input to search bar
16     search_button = driver.find_element_by_xpath('//*[@id="container"]/div/div[1]/div[1]/div[2]/div[2]/form/div/button')
17     search_button.click()
18 except NoSuchElementException:
19     driver.get(url)

```

```

1 links=[]
2 try:
3     for l in driver.find_elements_by_xpath('//a[@class="_1fQZEK"]'):
4         links.append(l.get_attribute('href'))
5         time.sleep(1)
6 except NoSuchElementException:
7     links.append("-")
8     time.sleep(1)
9
10 for page in range(0,10):
11     nxt_button=driver.find_elements_by_xpath("//a[@class='_1LKTO3']")#scraping the list of buttons from the page
12     time.sleep(2)
13     try:
14         driver.get(nxt_button[1].get_attribute('href'))#getting the link from the list for next page
15         time.sleep(2)
16     except:
17         driver.get(nxt_button[0].get_attribute('href'))
18         time.sleep(2)
19     try:
20         for l in driver.find_elements_by_xpath('//a[@class="_1fQZEK"]'):
21             links.append(l.get_attribute('href'))
22             time.sleep(1)
23     except NoSuchElementException:
24         links.append("-")
25         time.sleep(1)

```

```

1 len(links)

```

```

Rating=[]
Review=[]
Review_title=[]
for url in links:
    driver.get(url)
    time.sleep(2)
    try:
        all_reviews = driver.find_element_by_xpath('//div[@class="col JOpGWq"]/a')
        a = all_reviews.get_attribute('href')
        driver.get(a)
        driver.find_element_by_xpath('//*[@id="container"]/div/div[3]/div/div[1]/div[2]/div[1]/div/div[2]/div/div/div/div/se
        driver.find_element_by_xpath('//select[@class="_1EDlbo tVKh2S"]/option[@value="NEGATIVE_FIRST"]').click()
        driver.get(driver.current_url)

        # Extracting stars
        try:
            rating = driver.find_elements_by_xpath('//div[@class="col _2wzgFH K0kLPL"]/div[1]/div[1]')
            for r in rating:
                Rating.append(r.text)
            except NoSuchElementException:
                Rating.append('-')
            except StaleElementReferenceException:
                Rating.append('-')

            time.sleep(1)
            # Extracting review
            try:
                review = driver.find_elements_by_xpath('//div[@class="t-ZTKy"]/div/div')
                for r in review:
                    Review.append(r.text)
            except NoSuchElementException:
                Review.append('-')

            time.sleep(1)
            # Extracting review
            try:
                review_title = driver.find_elements_by_xpath('//p[@class="_2-N8zT"]')
                for r in review_title:
                    Review_title.append(r.text)
            except NoSuchElementException:
                Review_title.append('-')
        except TimeoutException:
            pass
        except NoSuchElementException:
            pass

```

Amazon

```

1 #webdriver
2 driver=webdriver.Chrome(r"C:\Users\SAGAR KADAM\Downloads\chromedriver_win32 (2)\chromedriver.exe")
3 time.sleep(3)

```

- TV

```

1 # Opening the homepage of Amazon.in
2 url = "https://www.amazon.in/"
3 driver.get(url)
4
5 time.sleep(2)
6 search_bar = driver.find_element_by_id("twotabsearchtextbox") # Locating searc_bar by id
7 search_bar.clear() # clearing search_bar
8 search_bar.send_keys('TV') # sending user input to search bar
9 search_button = driver.find_element_by_xpath('//div[@class="nav-search-submit nav-sprite"]/span/input') # Locating sea
10 search_button.click()
11 links=[]
12 try:
13     for l in driver.find_elements_by_xpath('//a[@class="a-link-normal a-text-normal"]'):
14         links.append(l.get_attribute('href'))
15         time.sleep(1)
16 except NoSuchElementException:
17     links.append("-")
18

```

```

1 len(links)

```

```

Rating=[]
Review=[]
Review_title=[]
for url in links:
    driver.get(url)
    time.sleep(2)
    try:
        all_reviews = driver.find_element_by_xpath('//*[@id="reviews-medley-footer"]/div[2]/a')
        a = all_reviews.get_attribute('href')
        driver.get(a)
        driver.find_element_by_xpath('//*[@id="a-autoid-6-announce"]').click()
        driver.find_element_by_xpath('//*[@id="star-count-dropdown_4"]').click()
        driver.get(driver.current_url)
        # Extracting stars
        try:
            rating = driver.find_elements_by_xpath('//*[@data-hook="review-star-rating"]/span')
            for r in rating:
                Rating.append(2)
        except NoSuchElementException:
            Rating.append('-')

        time.sleep(1)
        # Extracting review
        try:
            review = driver.find_elements_by_xpath('//*[@class="a-section celwidget"]/div[4]')
            for r in review:
                Review.append(r.text)
        except NoSuchElementException:
            Review.append('-')

        time.sleep(1)
        # Extracting review
        try:
            review_title = driver.find_elements_by_xpath('//*[@data-hook="review-title"]/span')
            for r in review_title:
                Review_title.append(r.text)
        except NoSuchElementException:
            Review_title.append('-')

    except NoSuchElementException:
        pass

```

```

1 # dictionary of lists
2 dict5 = {'Rating': total_R, 'Review': total_Re, 'Review_title': total_Rt}
3
4 df = pd.DataFrame(dict5)
5
6 # saving the dataframe
7 df.to_csv('Amazon.csv')

```

```

1 filename=['Nykaa.csv', 'Snapdeal.csv', 'Flipkart1.csv', 'Flipkart2.csv', 'Amazon.csv']

```

```

1 combined_csv = pd.concat([pd.read_csv(f) for f in filename])

```

```

1 combined_csv.to_excel( "Rating.xlsx", index=False, encoding='utf-8')

```

```

1 combined_csv.to_excel( "Rating_combined.xlsx", index=False, encoding='utf-8')

```


Firstly, we will start by importing required libraries and databases.

```
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
from wordcloud import WordCloud

from sklearn.linear_model import LogisticRegression
from sklearn.svm import SVC
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import AdaBoostClassifier, RandomForestClassifier

from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
from sklearn.model_selection import train_test_split
from sklearn.model_selection import cross_val_score

import joblib
import warnings
warnings.filterwarnings('ignore')
```

```
1 pd.set_option('display.max_colwidth',100 )
2 cd=pd.read_excel("Rating.xlsx")
3 data=pd.DataFrame(data=cd)
4 data
```

	Unnamed: 0	Rating	Review	Review_title
0	0	5star	Earlier i was using apple ipod but sound was not that much good and today i got this. Seriously ...	"Amazing sound "
1	1	5star	I mean for this price.... they're really nice one Good quality... Aesthetically pleasing and sou...	"Nice "
2	2	5star	perfect fitting nice sound quality awesome	"Amazing Product"
3	3	5star	Decent product. Worth giving a try!	"Good"
4	4	5star	Very good colour for girls . Nice voice i liked it. Price is good	"Very good product"
...
29016	1089	2	When I press the input mode change buttons on speaker or the remote there are only three options...	Sounds Quality is good but USB/SD mode not working
29017	1090	2	It's an amazing product but the one thing lack is bass control its poor. At this price range it ...	Ok ok it's sound is good but base quality very poor.
29018	1091	2	sound quality is ok bass is not good and if you are buying it for Bluetooth dont buy the worst B...	For wireless worst product
29019	1092	2	Worst ever music system just in a month it's 2 satellite audio Jack's got bad. Now i have to car...	Got bad in a month 😞. They are hundreds of complaints don't buy. And don't ignore complaints.
29020	1093	2	Firstly when I ordered this product I wanted to give my friend this as gift for his new shop. Th...	Low sound quality

29021 rows × 4 columns

```
1 data.columns
```

```
Index(['Unnamed: 0', 'Rating', 'Review', 'Review_title'], dtype='object')
```

Above is the list of columns in our database.

```
1 data.shape
```

```
(29021, 4)
```

Our database consists of 4 columns and 29021 rows.

```
1 data.dtypes
```

```
Unnamed: 0      int64
Rating          object
Review          object
Review_title    object
dtype: object
```

Our dataset consists of integer and object type of data.

```
1 #Unnamed: 0 is not required for further analysis hence we will drop the column
2 data.drop('Unnamed: 0',axis=1, inplace=True)
```

```
1 data.dropna(inplace=True)#dropped rows which consists of null values
```

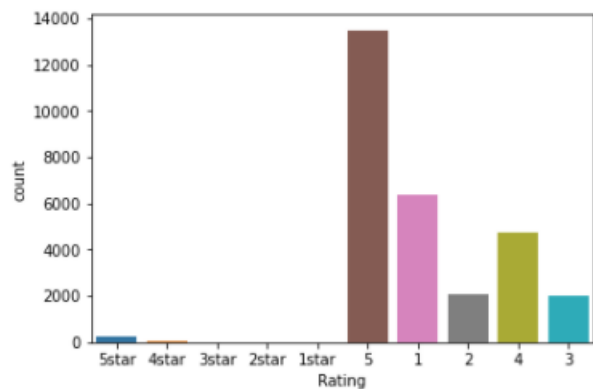
```
1 #lets check number of unique values in all columns
2 data.nunique()
```

```
Rating          10
Review         13350
Review_title    1672
dtype: int64
```

Our rating should only consist of 5 unique values i.e 1, 2, 3, 4, 5. let's have a look at rating column data with the help of countplot.

```
1 sns.countplot(data["Rating"])
```

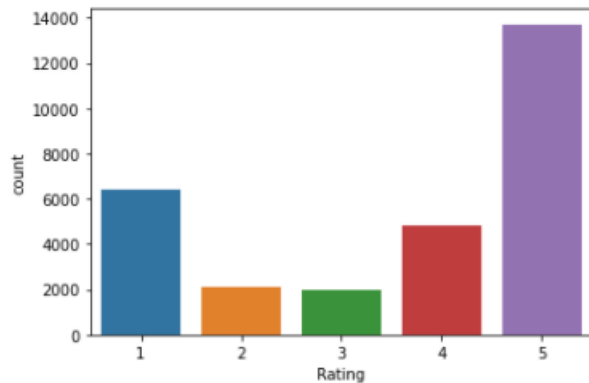
<matplotlib.axes._subplots.AxesSubplot at 0x23a33b1bf70>



```
1 data['Rating'] = data['Rating'].replace(['1star'],1, regex=True)#replacing 1star with 1
2 data['Rating'] = data['Rating'].replace(['2star'],2, regex=True)#replacing 2star with 2
3 data['Rating'] = data['Rating'].replace(['3star'],3, regex=True)#replacing 3star with 3
4 data['Rating'] = data['Rating'].replace(['4star'],4, regex=True)#replacing 4star with 4
5 data['Rating'] = data['Rating'].replace(['5star'],5, regex=True)#replacing 5star with 5
```

```
1 # reviewing rating column again
2 sns.countplot(data["Rating"])
```

<matplotlib.axes._subplots.AxesSubplot at 0x23a33d55580>



```
1 #lets look at counts of each values
2 data['Rating'].value_counts()
```

```
5    13709
1     6403
4     4802
2     2103
3     2003
Name: Rating, dtype: int64
```

```
1 data.shape
```

(29020, 3)

Let's do preprocessing of data with the help of NLP:

```
1 from nltk.corpus import stopwords
2 from nltk.tokenize import word_tokenize
3 import string
4 print("Stopwords :\n", stopwords.words("english"))
5 print("\n \n Punctuations :\n",string.punctuation)
```

Stopwords :

['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you're", "you've", "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'him', 'his', 'himself', 'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itself', 'they', 'them', 'their', 'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "that'll", 'these', 'those', 'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had', 'having', 'do', 'does', 'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as', 'until', 'while', 'of', 'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'through', 'during', 'before', 'after', 'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'over', 'under', 'again', 'further', 'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'any', 'both', 'each', 'few', 'more', 'most', 'other', 'some', 'such', 'no', 'nor', 'not', 'only', 'own', 'same', 'so', 'than', 'too', 'very', 's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "should've", 'now', 'd', 'll', 'm', 'o', 're', 've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn', "didn't", 'doesn', "doesn't", 'hadn', "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'ma', 'mightn', "mightn't", 'mustn', "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't", 'wasn', "wasn't", 'weren', "weren't", 'won', "won't", 'wouldn', "wouldn't"]

Punctuations :

!"#\$%&'()*+,-./:;<=>?@[\\]^_`{|}~

```
1 punct=string.punctuation
```

```
1 #function to remove punctuation
2 def remove_punctuation(txt):
3     txt_nopunct="".join([c for c in txt if c not in punct])
4     return txt_nopunct
```

```
1 #function to remove stopwords
2 def rem_stopworrrds(text):
3     stop_words=set(stopwords.words("english"))
4     word_tokens=word_tokenize(text)
5     filtered_text=" ".join([word for word in word_tokens if word not in stop_words])
6     return filtered_text
```

```
1 #Let's preprocss all 3 columns
2 data['Review'] = data['Review'].apply(lambda x: x.lower())# Lets change the review into lowercase
3 data['Review_title'] = data['Review_title'].apply(lambda x: x.lower())# Lets change the review_title into lowercase
4
5 data['Review'] = data['Review'].apply(lambda x: remove_punctuation(x))
6 data['Review'] = data['Review'].apply(lambda x: rem_stopworrrds(x))
7
8 data['Review_title'] = data['Review_title'].apply(lambda x: remove_punctuation(x))
9 data['Review_title'] = data['Review_title'].apply(lambda x: remove_punctuation(x))
```

```
1 data['Review']
```

```
0      earlier using apple ipod sound much good today got seriously ipod amazing sound quality bass eve...
1          mean price theyre really nice one good quality aesthetically pleasing sound quality good
2                                perfect fitting nice sound quality awesome
3                                decent product worth giving try
4                                good colour girls nice voice liked price good
5                                ...
29016  press input mode change buttons speaker remote three options aux bluetooth fm option usb sd card...
29017  amazing product one thing lack bass control poor price range could given better features atleast...
29018  sound quality ok bass good buying bluetooth dont buy worst bluetooth quality range buy portable ...
29019  worst ever music system month 2 satellite audio jacks got bad carry service centres dont know ki...
29020  firstly ordered product wanted give friend gift new shop product delivered time woofer working r...
Name: Review, Length: 29020, dtype: object
```

```
1 import nltk
2 wn=nltk.WordNetLemmatizer()
```

```
1 #Lemmatization function
2 def lemmatization(txt):
3     text="".join(wn.lemmatize(c) for c in txt)
4     return text
```

```
1 #Lemmatized the reviews
2 data['Review'] = data['Review'].apply(lambda x: lemmatization(x))
```

```
1 data['Review']
```

```
0      earlier using apple ipod sound much good today got seriously ipod amazing sound quality bass eve...
1          mean price theyre really nice one good quality aesthetically pleasing sound quality good
2                                perfect fitting nice sound quality awesome
3                                decent product worth giving try
4                                good colour girls nice voice liked price good
5                                ...
29016  press input mode change buttons speaker remote three options aux bluetooth fm option usb sd card...
29017  amazing product one thing lack bass control poor price range could given better features atleast...
29018  sound quality ok bass good buying bluetooth dont buy worst bluetooth quality range buy portable ...
29019  worst ever music system month 2 satellite audio jacks got bad carry service centres dont know ki...
29020  firstly ordered product wanted give friend gift new shop product delivered time woofer working r...
Name: Review, Length: 29020, dtype: object
```

Let's join all words in review column in a variable text:

```
1 text = " ".join(review for review in data['Review'])
2 print ("There are {} words in the combination of all review.".format(len(text)))
```

There are 1862925 words in the combination of all review.

```
1 stopwords_en = set(stopwords.words('english'))
```

Wordcloud

Word cloud of word in text variable:

```
1 # Create stopwords list:
2 stopwords = set(stopwords_en)
3
4 # Generate a word cloud image
5 wordcloud = WordCloud(stopwords=stopwords, background_color="white").generate(text)
6
7 # Display the generated image:
8 # the matplotlib way:
9 plt.figure(figsize=(10,8),facecolor='r')
10 plt.imshow(wordcloud, interpolation='bilinear')
11 plt.axis("off")
12 plt.show()
```



Above cloud gives us the most common words in entire reviews in bigfont size, word with smaller size is less used or less common.

```
1 #Getting sense of loud wrds in spam
2 from wordcloud import WordCloud
3
4
5 view = data['Review'][data['Rating']==1]
6
7 spam_cloud = WordCloud(width=700,height=500,background_color='white',max_words=20).generate(' '.join(view))
8
9 plt.figure(figsize=(10,8),facecolor='r')
10 plt.imshow(spam_cloud)
11 plt.axis('off')
12 plt.tight_layout(pad=0)
13 plt.show()
```




```
1 from imblearn import under_sampling
2 from collections import Counter
```

```
1 from imblearn.under_sampling import RandomUnderSampler
2 |
3 rus=RandomUnderSampler(random_state=0)
4 x_resample,y_resampled=rus.fit_resample(x,y)
5 print(sorted(Counter(y_resampled)),y_resampled.shape)
```

```
[1, 2, 3, 4, 5] (10015,)
```

```
1 y_resampled.value_counts()
```

```
5    2003
4    2003
3    2003
2    2003
1    2003
Name: Rating, dtype: int64
```

```
1 x_resample.shape
```

```
(10015, 13343)
```

```
1 y_resampled.shape
```

```
(10015,)
```

```
1 x=x_resample
2 y=y_resampled
```

```
1 x_train,x_test, y_train,y_test=train_test_split(x,y,test_size=0.2, random_state=0)
```

```
1 x_train.shape, y_train.shape, x_test.shape, y_test.shape
```

```
((8012, 13343), (8012,), (2003, 13343), (2003,))
```

Let's find best random state:

```
1 maxAccu=0
2 maxRS=0
3 for i in range(1,200):
4     x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=.20,random_state=i)
5     LR=LogisticRegression()
6     LR.fit(x_train,y_train)
7     pred=LR.predict(x_test)
8     acc=accuracy_score(y_test,pred)
9     if acc>maxAccu:
10         maxAccu=acc
11         maxRS=i
12 print("Best accuracy is ",maxAccu, " on Random State ",maxRS)
```

```
Best accuracy is  0.5831253120319521  on Random State  21
```


Model/s Development and Evaluation

```
1 x_train,x_test, y_train,y_test=train_test_split(x,y,test_size=0.2, random_state=maxRS)

1 model=[LogisticRegression(),DecisionTreeClassifier(),SVC(),AdaBoostClassifier(),RandomForestClassifier()]
2
3 for m in model:
4     m.fit(x_train,y_train)
5     #m.score(x_train,y_train)
6     pred=m.predict(x_test)
7     acc=accuracy_score(y_test,pred)
8     print('Accuracy Score of',m,'is:', acc)
9     print(confusion_matrix(y_test,pred))
10    print(classification_report(y_test,pred))
11    print('\n')
```

Accuracy Score of LogisticRegression() is: 0.5831253120319521

```
[[273  61  31  15   8]
 [ 67 253  46  13   2]
 [ 62  53 186  72  43]
 [ 21  23  47 231  88]
 [ 14  20  31 118 225]]
      precision    recall  f1-score   support

     1         0.62      0.70      0.66         388
     2         0.62      0.66      0.64         381
     3         0.55      0.45      0.49         416
     4         0.51      0.56      0.54         410
     5         0.61      0.55      0.58         408

 accuracy          0.58
macro avg          0.58      0.59      0.58         2003
weighted avg          0.58      0.58      0.58         2003
```

Accuracy Score of DecisionTreeClassifier() is: 0.5621567648527209

```
[[243  60  39  24  22]
 [ 72 201  50  36  22]
 [ 44  32 236  71  33]
 [ 23  21  43 258  65]
 [ 33  26  38 123 188]]
      precision    recall  f1-score   support

     1         0.59      0.63      0.61         388
     2         0.59      0.53      0.56         381
     3         0.58      0.57      0.57         416
     4         0.50      0.63      0.56         410
     5         0.57      0.46      0.51         408

 accuracy          0.56
macro avg          0.57      0.56      0.56         2003
weighted avg          0.57      0.56      0.56         2003
```

Accuracy Score of SVC() is: 0.6380429355966051

```
[[302  44  17  18   7]
 [ 66 260  39  14   2]
 [ 53  39 217  72  35]
 [ 15  15  38 262  80]
 [ 12  12  27 120 237]]
      precision    recall  f1-score   support

     1         0.67      0.78      0.72         388
     2         0.70      0.68      0.69         381
     3         0.64      0.52      0.58         416
     4         0.54      0.64      0.58         410
     5         0.66      0.58      0.62         408

 accuracy          0.64
macro avg          0.64      0.64      0.64         2003
weighted avg          0.64      0.64      0.64         2003
```

Accuracy Score of AdaBoostClassifier() is: 0.45481777333999

```
[[228  83  51  12  14]
 [ 83 192  76  19  11]
 [ 69  94 118  73  62]
 [ 13  29  68 132 168]
 [  9  24  47  87 241]]
      precision    recall  f1-score   support

     1         0.57      0.59      0.58         388
     2         0.45      0.50      0.48         381
     3         0.33      0.28      0.30         416
     4         0.41      0.32      0.36         410
     5         0.49      0.59      0.53         408

 accuracy
macro avg         0.45      0.46      0.45        2003
weighted avg         0.45      0.45      0.45        2003
```

Accuracy Score of RandomForestClassifier() is: 0.6255616575137294

```
[[293  55  15  17   8]
 [ 73 259  30  11   8]
 [ 47  51 221  66  31]
 [ 28  18  31 255  78]
 [ 23  26  26 108 225]]
      precision    recall  f1-score   support

     1         0.63      0.76      0.69         388
     2         0.63      0.68      0.66         381
     3         0.68      0.53      0.60         416
     4         0.56      0.62      0.59         410
     5         0.64      0.55      0.59         408

 accuracy
macro avg         0.63      0.63      0.62        2003
weighted avg         0.63      0.63      0.62        2003
```

Hyper parameter tuning:

```
1 from sklearn.model_selection import RandomizedSearchCV
2 #creating parameter list to pass in RandomizedSearchCV
```

- LogisticRegression

```
1 #DecisionTreeClassifier parameters
2 parameters1={'penalty':['l2','l1','elasticnet', 'none'], 'dual':[True,False], 'tol':[0.0001], 'C':[1.0],
3              'fit_intercept':[True,False], 'intercept_scaling':[1], 'class_weight':[None], 'random_state':range(0,20),
4              'solver':['newton-cg', 'lbfgs', 'liblinear', 'sag', 'saga'], 'max_iter':[100],
5              'multi_class':['auto', 'ovr', 'multinomial'], 'verbose':[0], 'warm_start':[True,False],
6              'n_jobs':[None], 'l1_ratio':[None]}
```

```
1 RSV1=RandomizedSearchCV(LogisticRegression(),parameters1,cv=5)
```

```
1 RSV1.fit(x_train,y_train)
```

```
RandomizedSearchCV(cv=5, estimator=LogisticRegression(),
                  param_distributions={'C': [1.0], 'class_weight': [None],
                                      'dual': [True, False],
                                      'fit_intercept': [True, False],
                                      'intercept_scaling': [1],
                                      'l1_ratio': [None], 'max_iter': [100],
                                      'multi_class': ['auto', 'ovr',
                                                    'multinomial'],
                                      'n_jobs': [None],
                                      'penalty': ['l2', 'l1', 'elasticnet',
                                                  'none'],
                                      'random_state': range(0, 20),
                                      'solver': ['newton-cg', 'lbfgs',
                                                  'liblinear', 'sag', 'saga'],
                                      'tol': [0.0001], 'verbose': [0],
                                      'warm_start': [True, False]})
```

```
1 RSV1.best_params_
```

```
{'warm_start': True,
 'verbose': 0,
 'tol': 0.0001,
 'solver': 'sag',
 'random_state': 19,
 'penalty': 'none',
 'n_jobs': None,
 'multi_class': 'ovr',
 'max_iter': 100,
 'l1_ratio': None,
 'intercept_scaling': 1,
 'fit_intercept': True,
 'dual': False,
 'class_weight': None,
 'C': 1.0}
```

```
1 RSV_pred1=RSV1.best_estimator_.predict(x_test)
```

```
1 RSV_pred1
```

```
array([5, 2, 5, ..., 1, 3, 4], dtype=int64)
```

```
1 RSV1.score(x_train,y_train)
```

```
0.8578382426360459
```

Hyper parameter tuning with LogisticRegression gives us 85.78 % accuracy.

- DecisionTreeClassifier

```
1 #DecisionTreeClassifier parameters
2 parameters2={'criterion':['gini', 'entropy'],'splitter':['best', 'random'],
3              'max_depth':[None],'random_state':range(0,20),
4              'max_leaf_nodes':[None],'class_weight':[None], 'min_samples_split':[1,2,3],
5              'min_samples_leaf':[1,2,3],'max_features':[None],'min_impurity_decrease':[0.0,1.0,2.0],
6              'min_weight_fraction_leaf':[0.0,1.0,2.0]}
```

```
1 RSV2=RandomizedSearchCV(DecisionTreeClassifier(),parameters2,cv=5)
```

```
1 RSV2.fit(x_train,y_train)
```

```
RandomizedSearchCV(cv=5, estimator=DecisionTreeClassifier(),
                  param_distributions={'class_weight': [None],
                                      'criterion': ['gini', 'entropy'],
                                      'max_depth': [None],
                                      'max_features': [None],
                                      'max_leaf_nodes': [None],
                                      'min_impurity_decrease': [0.0, 1.0,
                                                                2.0],
                                      'min_samples_leaf': [1, 2, 3],
                                      'min_samples_split': [1, 2, 3],
                                      'min_weight_fraction_leaf': [0.0, 1.0,
                                                                    2.0],
                                      'random_state': range(0, 20),
                                      'splitter': ['best', 'random']})
```

```
1 RSV2.best_params_
```

```
{'splitter': 'best',
 'random_state': 14,
 'min_weight_fraction_leaf': 0.0,
 'min_samples_split': 2,
 'min_samples_leaf': 2,
 'min_impurity_decrease': 0.0,
 'max_leaf_nodes': None,
 'max_features': None,
 'max_depth': None,
 'criterion': 'entropy',
 'class_weight': None}
```

```
1 RSV_pred2=RSV2.best_estimator_.predict(x_test)
```

```

1 RSV_pred2
array([5, 1, 5, ..., 1, 5, 3], dtype=int64)

1 RSV2.score(x_train,y_train)
0.8358711932101848

```

Hyper parameter tuning with DecisionTreeClassifier gives us 83.59 % accuracy.

- SCV

```

1 #SVC parameters
2 parameters3={'C':[1, 10, 100, 1000], 'kernel':['linear', 'poly', 'rbf', 'sigmoid'],
3             'degree':[3], 'gamma':['scale', 'auto'], 'coef0':[0.0], 'shrinking':[True, False],
4             'probability':[True, False], 'tol':[0.001], 'cache_size':[200],
5             'class_weight':[None], 'decision_function_shape':['ovr'],
6             'break_ties':[True, False], 'random_state':range(0,20)}

1 RSV3=RandomizedSearchCV(SVC(),parameters3,cv=5)

1 RSV3.fit(x_train,y_train)

RandomizedSearchCV(cv=5, estimator=SVC(),
                  param_distributions={'C': [1, 10, 100, 1000],
                                      'break_ties': [True, False],
                                      'cache_size': [200],
                                      'class_weight': [None], 'coef0': [0.0],
                                      'decision_function_shape': ['ovr'],
                                      'degree': [3],
                                      'gamma': ['scale', 'auto'],
                                      'kernel': ['linear', 'poly', 'rbf',
                                                'sigmoid'],
                                      'probability': [True, False],
                                      'random_state': range(0, 20),
                                      'shrinking': [True, False],
                                      'tol': [0.001]}))

1 RSV3.best_params_
{'tol': 0.001,
 'shrinking': True,
 'random_state': 6,
 'probability': False,
 'kernel': 'poly',
 'gamma': 'scale',
 'degree': 3,
 'decision_function_shape': 'ovr',
 'coef0': 0.0,
 'class_weight': None,
 'cache_size': 200,
 'break_ties': False,
 'C': 1}

1 RSV_pred3=RSV3.best_estimator_.predict(x_test)

1 RSV_pred3
array([5, 2, 5, ..., 1, 4, 4], dtype=int64)

1 RSV3.score(x_train,y_train)
0.858337493759361

```

Hyper parameter tuning with SVC gives us 85.83 % accuracy.

- AdaBoostClassifier

```
1 parameters4={'base_estimator':[None], 'n_estimators':[20,40, 50,75,100],
2             'learning_rate':[0.1,0.01,1.0,2.0], 'algorithm':['SAMME.R', 'SAMME'],
3             'random_state':range(0,20)}
4
```

```
1 RSV4=RandomizedSearchCV(AdaBoostClassifier(),parameters4,cv=5)
```

```
1 RSV4.fit(x_train,y_train)
```

```
RandomizedSearchCV(cv=5, estimator=AdaBoostClassifier(),
                  param_distributions={'algorithm': ['SAMME.R', 'SAMME'],
                                      'base_estimator': [None],
                                      'learning_rate': [0.1, 0.01, 1.0, 2.0],
                                      'n_estimators': [20, 40, 50, 75, 100],
                                      'random_state': range(0, 20)})
```

```
1 RSV4.best_params_
```

```
{'random_state': 0,
 'n_estimators': 100,
 'learning_rate': 2.0,
 'base_estimator': None,
 'algorithm': 'SAMME'}
```

```
1 RSV_pred4=RSV4.best_estimator_.predict(x_test)
```

```
1 RSV_pred4
```

```
array([4, 2, 5, ..., 2, 3, 4], dtype=int64)
```

```
1 RSV4.score(x_train,y_train)
```

```
0.4162506240639041
```

Hyper parameter tuning with adaboost classifier gives us 41.63 % accuracy.

- RandomForestClassifier

```
: 1 parameters5={'n_estimators':[100], 'criterion':['gini', 'entropy'], 'max_depth':[None],
2             'min_samples_split':[2], 'min_samples_leaf':[1], 'min_weight_fraction_leaf':[0.0],
3             'max_features':['auto'], 'max_leaf_nodes':[None], 'min_impurity_decrease':[0.0],
4             'min_impurity_split':[None], 'bootstrap':[True, False], 'oob_score':[True, False], 'n_jobs':[None],
5             'random_state':range(0,20), 'verbose':[0], 'warm_start':[True, False], 'class_weight':[None],
6             'ccp_alpha':[0.0], 'max_samples':[None]}
```

```
: 1 RSV5=RandomizedSearchCV(RandomForestClassifier(),parameters5,cv=5)
```

```
: 1 RSV5.fit(x_train,y_train)
```

```
RandomizedSearchCV(cv=5, estimator=RandomForestClassifier(),
                  param_distributions={'bootstrap': [True, False],
                                      'ccp_alpha': [0.0],
                                      'class_weight': [None],
                                      'criterion': ['gini', 'entropy'],
                                      'max_depth': [None],
                                      'max_features': ['auto'],
                                      'max_leaf_nodes': [None],
                                      'max_samples': [None],
                                      'min_impurity_decrease': [0.0],
                                      'min_impurity_split': [None],
                                      'min_samples_leaf': [1],
                                      'min_samples_split': [2],
                                      'min_weight_fraction_leaf': [0.0],
                                      'n_estimators': [100], 'n_jobs': [None],
                                      'oob_score': [True, False],
                                      'random_state': range(0, 20),
                                      'verbose': [0],
                                      'warm_start': [True, False]})
```

```
: 1 RSV5.best_params_
```

```
{'warm_start': False,
 'verbose': 0,
 'random_state': 12,
 'oob_score': False,
 'n_jobs': None,
 'n_estimators': 100,
 'min_weight_fraction_leaf': 0.0,
 'min_samples_split': 2,
 'min_samples_leaf': 1,
 'min_impurity_split': None,
 'min_impurity_decrease': 0.0,
 'max_samples': None,
 'max_leaf_nodes': None,
 'max_features': 'auto',
 'max_depth': None,
 'criterion': 'gini',
 'class_weight': None,
 'ccp_alpha': 0.0,
 'bootstrap': True}
```

```
1 RSV_pred5=RSV5.best_estimator_.predict(x_test)
```

```
1 RSV_pred5
```

```
array([5, 2, 5, ..., 1, 1, 4], dtype=int64)
```

```
1 RSV5.score(x_train,y_train)
```

```
0.8775586620069895
```

Hyper parameter tuning with RandomForestClassifier gives us 87.76% accuracy.

Considering RandomForestClassifier with RandomizedSearchCV as final model for saving model.

Saving model:

```
1 import joblib
2 joblib.dump(RSV5,"RSVRPP.obj")
```

```
['RSVRPP.obj']
```

```
1 RSVfile=joblib.load("RSVRPP.obj")
2 RSVfile.predict(x_test)
```

```
array([5, 2, 5, ..., 1, 1, 4], dtype=int64)
```

CONCLUSION

In the Ratings Prediction Project, I have extracted ratings and reviews from different websites then saved the extracted csv file into one excel sheet. Our dataset mainly consists of 3 columns and 29021 rows. Columns are Rating, Review and Review title though we have one more column named Unnamed: 0 however that is not required for further analysis.

Then I did some preprocessing like dropping null values and Our rating should only consist of 5 unique values i.e 1, 2, 3, 4, 5 and the rating column actually consists of 10 unique values hence I used the replace method. Then with the help of NLP I have removed punctuations, stopwords from Review and review title column. Then used a lemmatization method to change words with similar meaning to one common word. With the help of wordcloud I performed some visualization.

I have then used the TFidf method to vectorise and transform data. Values in our target column have different value counts hence with the use of random under sampler data of target column came down to equal value count. Used 5 methods for model building then with the help of RandomizedSearchCV I have tried to improve accuracy. Finally, I decided to go ahead with the Random Forest Classifier and saved the model.