



# **Ratings Prediction Project**

Submitted By  
Mr. Shubham J. Chougule

## **ACKNOWLEDGMENT**

It is my sensual gratification to present this report on RATINGS PREDICTION project which is a NLP project. Working on this project was a good experience that has given me a very informative knowledge.

I would like to express my sincere thanks to Ms. Khushboo Garg for a regular follow up and valuable guidance provided throughout.

And I am also thankful to FlipRobo Technologies Bangalore for their guidance and constant supervision as well as for providing necessary information regarding the project and also for their support in completing the project.

# Introduction:

## Business Problem Framing:

The rise in E-commerce has brought a significant rise in the importance of customer reviews. There are hundreds of review sites online and massive amounts of reviews for every product. Customers have changed their way of shopping and according to a recent survey, 70 percent of customers say that they use rating filters to filter out low rated items in their searches.

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, companies like Google, Amazon.

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. The second one is based on recommender systems, specifically on collaborative filtering, and focuses on the reviewer's point of view.

## Conceptual background of domain problem

Rating prediction is a well-known recommendation task aiming to predict a user's rating for items which were not yet rated. Predictions are computed from users' explicit feedback, i.e. their ratings provided on some items in the past. Another type of feedback are user reviews provided on items which implicitly express users' opinions on items. Recent studies indicate that opinions inferred from users' reviews on items are strong predictors of user's implicit feedback or even ratings and thus, should be utilized in computation. As far as we know, all the recent works on recommendation techniques utilizing opinions inferred from users' reviews are either focused on the item recommendation task or use only the opinion information, completely leaving users' ratings out of consideration. The approach proposed in this paper is filling this gap, providing a simple, personalized and scalable rating prediction framework utilizing both ratings provided by users and opinions inferred from their reviews.

Experimental results provided on dataset containing user ratings and reviews from the real-world Amazon and Flipkart Product Review Data show the effectiveness of the proposed framework.

## Analytical Problem Framing:

### Mathematical/Analytical modeling of the problem

As per the client's requirement for this rating prediction project I have scraped reviews and ratings from well-known e-commerce sites. This is then saved into .csv format. Also, I have shared the script for web scraping into the GitHub repository.

Then loaded this data into a data frame and did some of the important natural language processing steps and gone through several EDA steps to analyze the data. After all the necessary steps I have build a NLP ML model to predict the ratings.

#### Data Loading

```
1 #loading the data into a dataframe
2 df = pd.read_csv("/Users/shubh/Desktop/Rating-Prediction-Project-main copy/Rating.csv")
3 df
```

[5]

...	Unnamed: 0	Review_title	Reiew_text	Ratings
	0	Suitable for School kids	\n If you are a College student or a professi...	2.0 out of 5 stars
	1	Misrepresentation on MS Office 2019 license - ...	\n Update after one month usage - MS Office 2...	2.0 out of 5 stars
	2	The sold me renewed laptop	\n It's look like renewed laptop because lapt...	2.0 out of 5 stars
	3	Amazon dupes with specification/ battery sucks	\n &nbsp;I had seen the specifications and bo...	2.0 out of 5 stars
	4	Display back light issue	\n Display gone with 2 months.. But anyway th...	2.0 out of 5 stars
	...	...	...	...
	77545	Nice product	good product	4
	77546	Awesome	Very good as expected and happy with the purchase	5
	77547	Awesome	I love it! No complaint!	5
	77548	Nice product	good product	4
	77549	Awesome	Very good as expected and happy with the purchase	5

77550 rows x 4 columns

```
1 #info
2 df.info()
```

[10]

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 77550 entries, 0 to 77549
Data columns (total 3 columns):
#   Column      Non-Null Count  Dtype
---  -
0   Review_title 68521 non-null object
1   Reiew_text   69465 non-null object
2   Ratings      68523 non-null object
dtypes: object(3)
memory usage: 1.8+ MB
```

Looking at above both figures we can see that our data set contains 77550 different rows and 4 columns among which I have removed unwanted column(Unnamed:0). And for this project Ratings is our target column. There are some missing values in our dataset which have been removed from the dataset.

## Data Processing:

At first I have joined both columns Review\_title and Review\_text into a new column as Review.

### Combine Review\_title and Reiew\_text into one

```
1 #joining Review text and title
2 df['Review'] = df['Review_title'].map(str)+' '+df['Reiew_text']
```

[19]

Then all the entries from Ratings columns have been converted to respective integer values

```
1 df['Ratings'] = df['Ratings'].replace('1.0 out of 5 stars',1)
2 df['Ratings'] = df['Ratings'].replace('2.0 out of 5 stars',2)
3 df['Ratings'] = df['Ratings'].replace('3.0 out of 5 stars',3)
4 df['Ratings'] = df['Ratings'].replace('4.0 out of 5 stars',4)
5 df['Ratings'] = df['Ratings'].replace('5.0 out of 5 stars',5)
6 df['Ratings'] = df['Ratings'].astype('int')
```

[17]

## Text processing

### Text Processing

```
1 #Here I am defining a function to replace some of the contracted words to their full form and removing urls and some unwanted text
2 def decontracted(text):
3     text = re.sub(r"won't", "will not", text)
4     text = re.sub(r"don't", "do not", text)
5     text = re.sub(r"can't", "can not", text)
6     text = re.sub(r"im ", "i am", text)
7     text = re.sub(r"yo ", "you ", text)
8     text = re.sub(r"doesn't", "does not", text)
9     text = re.sub(r"n't", " not", text)
10    text = re.sub(r"\re", " are", text)
11    text = re.sub(r"\s", " is", text)
12    text = re.sub(r"\d", " would", text)
13    text = re.sub(r"\ll", " will", text)
14    text = re.sub(r"\t", " not", text)
15    text = re.sub(r"\ve", " have", text)
16    text = re.sub(r"\m", " am", text)
17    text = re.sub(r"<br>", " ", text)
18    text = re.sub(r'http\S+', '', text) #removing urls
19    return text
```

[24]

For text processing I have defined a function to replace some of the words with proper words. All text is converted to lowercase and removed different punctuations from the text of Review column.

**Lemmatization:** Lemmatization is the process of grouping together the different inflected forms of a word so they can be analyzed as a single item. Lemmatization is similar to stemming but it brings context to the words.

#### Lemmatization

```
1 #Defining function to convert nltk tag to wordnet tags
2 def nltk_tag_to_wordnet_tag(nltk_tag):
3     if nltk_tag.startswith('J'):
4         return wordnet.ADJ
5     elif nltk_tag.startswith('V'):
6         return wordnet.VERB
7     elif nltk_tag.startswith('N'):
8         return wordnet.NOUN
9     elif nltk_tag.startswith('R'):
10        return wordnet.ADV
11    else:
12        return None
```

[31]

```
1 #defining function to lemmatize our text
2 def lemmatize_sentence(sentence):
3     #tokenize the sentence & find the pos tag
4     nltk_tagged = nltk.pos_tag(nltk.word_tokenize(sentence))
5     #tuple of (token, wordnet_tag)
6     wordnet_tagged = map(lambda x : (x[0], nltk_tag_to_wordnet_tag(x[1])), nltk_tagged)
7     lemmatize_sentence = []
8     for word, tag in wordnet_tagged:
9         if tag is None:
10            lemmatize_sentence.append(word)
11        else:
12            lemmatize_sentence.append(lemmatizer.lemmatize(word,tag))
13    return " ".join(lemmatize_sentence)
```

[32]

For lemmatizing the text I have defined these two functions first will give the wordnet tag for the nltk\_tagged word then with respect to this wordnet tag lemmatization of each word is done.

## Text Normalization – Standardization

### Text Normalization - Standardization

```
1 #Noise removal
2 def scrub_words(text):
3     #remove html markup
4     text = re.sub("<.*?>", "", text)
5     #remove non-ascii and digits
6     text = re.sub("[^\w]", " ", text)
7     text = re.sub("[\d]", "", text)
8     #remove white space
9     text = text.strip()
10    return text
```

[36]

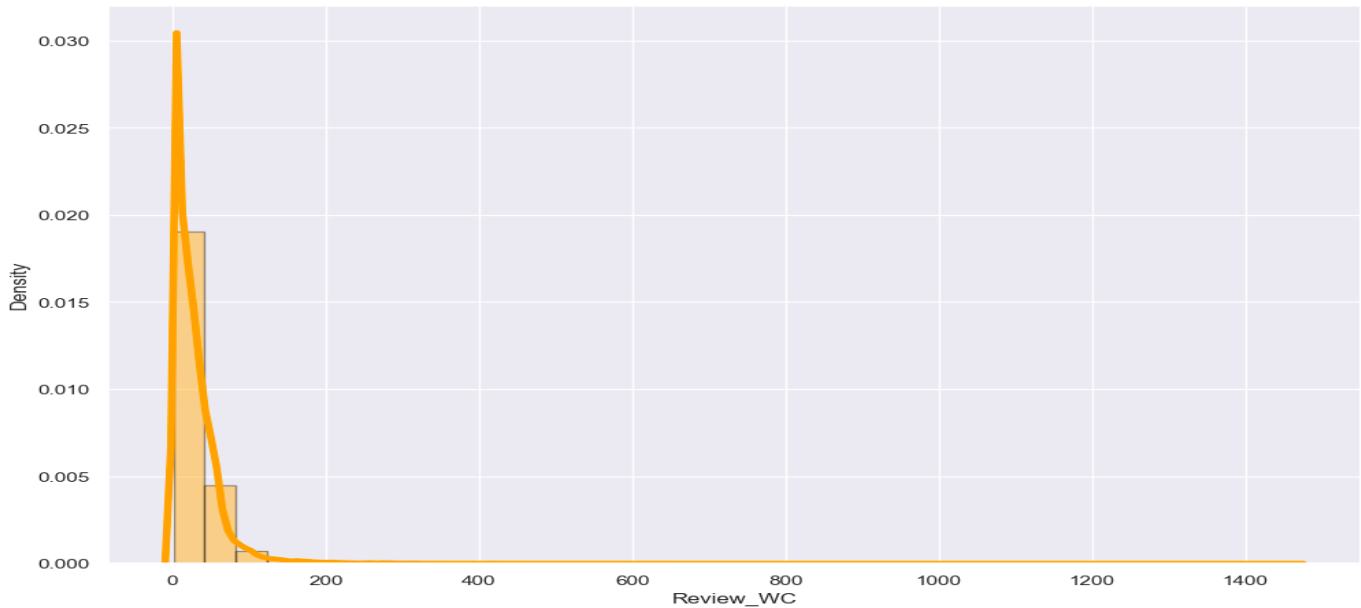
```
1 df['Review'] = df['Review'].apply(lambda x : scrub_words(x))
```

[37]

Finally for standardizing our text and removing numbers from it I have defined a function as scrub\_words as shown in above code and applied to the review column.

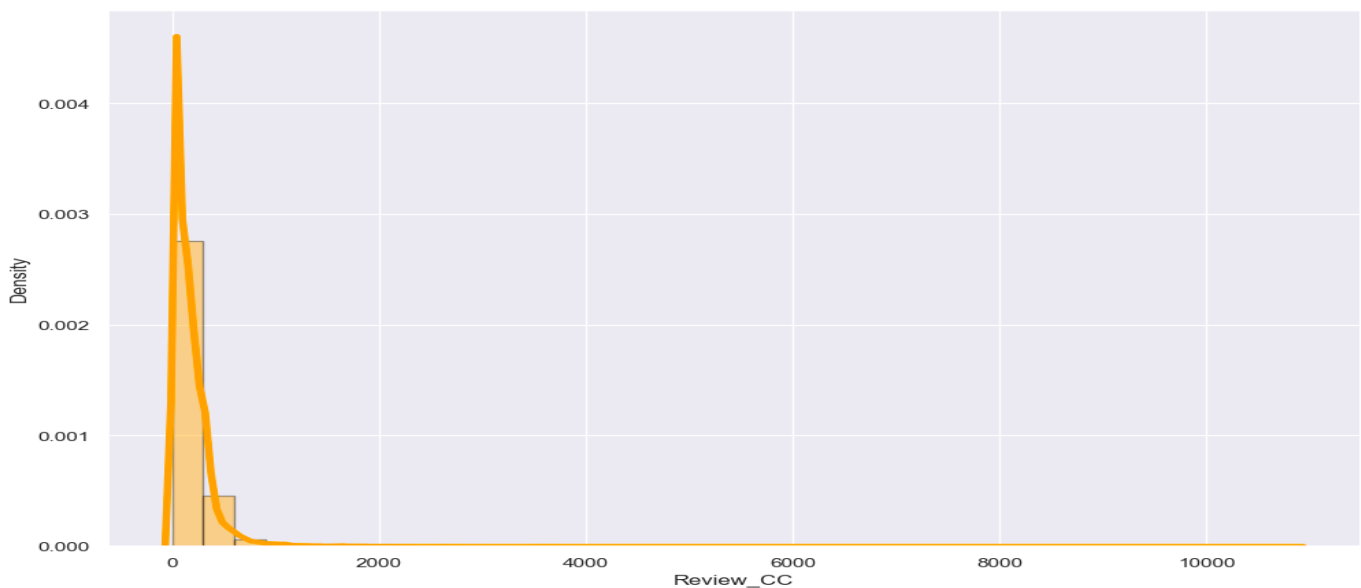
## Exploratory Data Analysis:

### Word\_count of review



Above figure shows the number of words from each review text. Looking at this histogram we can conclude that most of the review text is in the range of 0 to 200 of words. Rest reviews can be considered as outliers in our data.

### Character count of review



The plot for character count is almost similar to the plot of word count. We can see that most of the reviews are in the range of 0 to 1500 numbers of characters.

Looking at these plots I have decided to remove the data with too long reviews by considering them as outliers.

### Removing Outliers

As we know that some of the review are too lengthy I am removing those reviews from the data as outliers using `z_score` method.

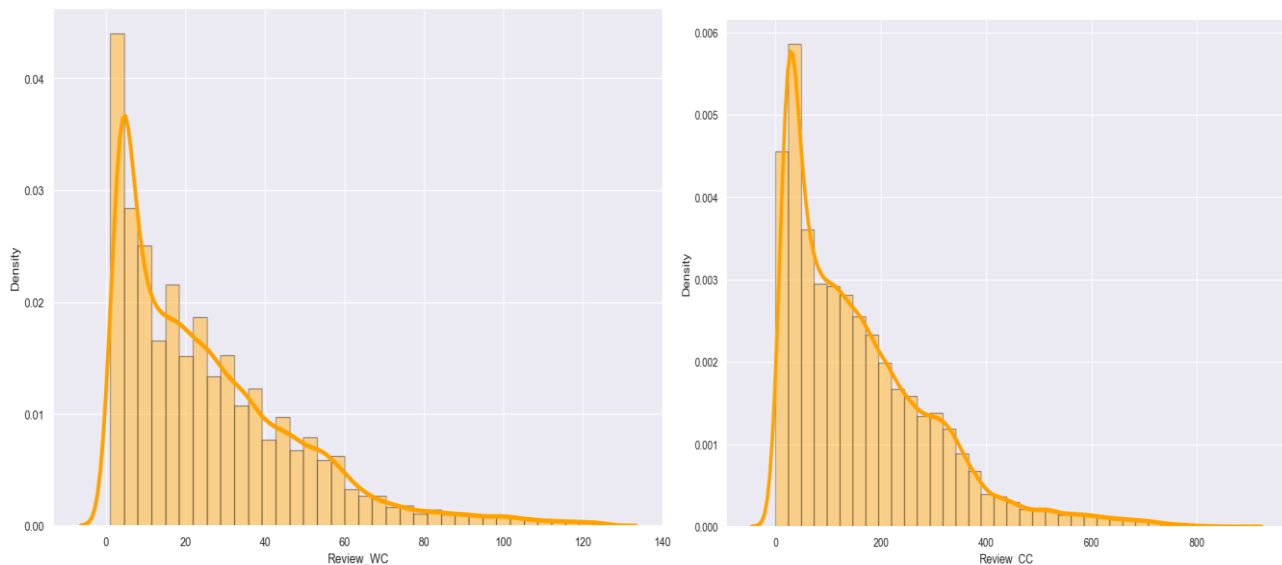
```
1 #apply zscore to remove outliers
2 from scipy import stats
3 from scipy.stats import zscore
4 z_score = zscore(df[['Review_WC']])
5 abs_z_score = np.abs(z_score)
6 filtering_entry = (abs_z_score < 3).all(axis = 1)
7 df = df[filtering_entry]
8 df.shape
```

[44]

... (67260, 6)

And by removing these outliers I am not losing much of the data so it is good to remove those entries for better results.

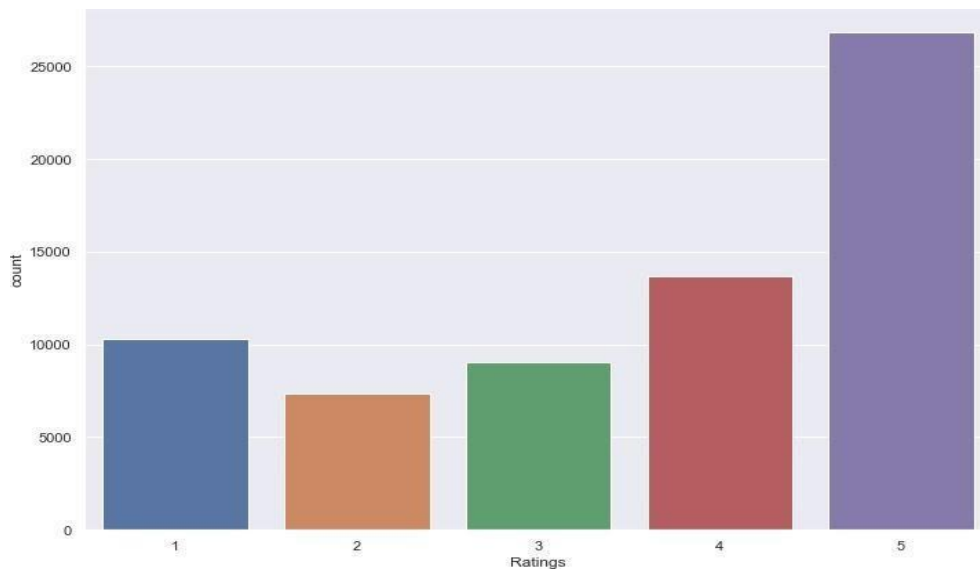
### Plotting histograms for word count and character counts again after removing outliers





After plotting histograms for word counts and character counts after removing outliers we can see now we are with good range of number of words and characters.

### Ratings (Target Variable):

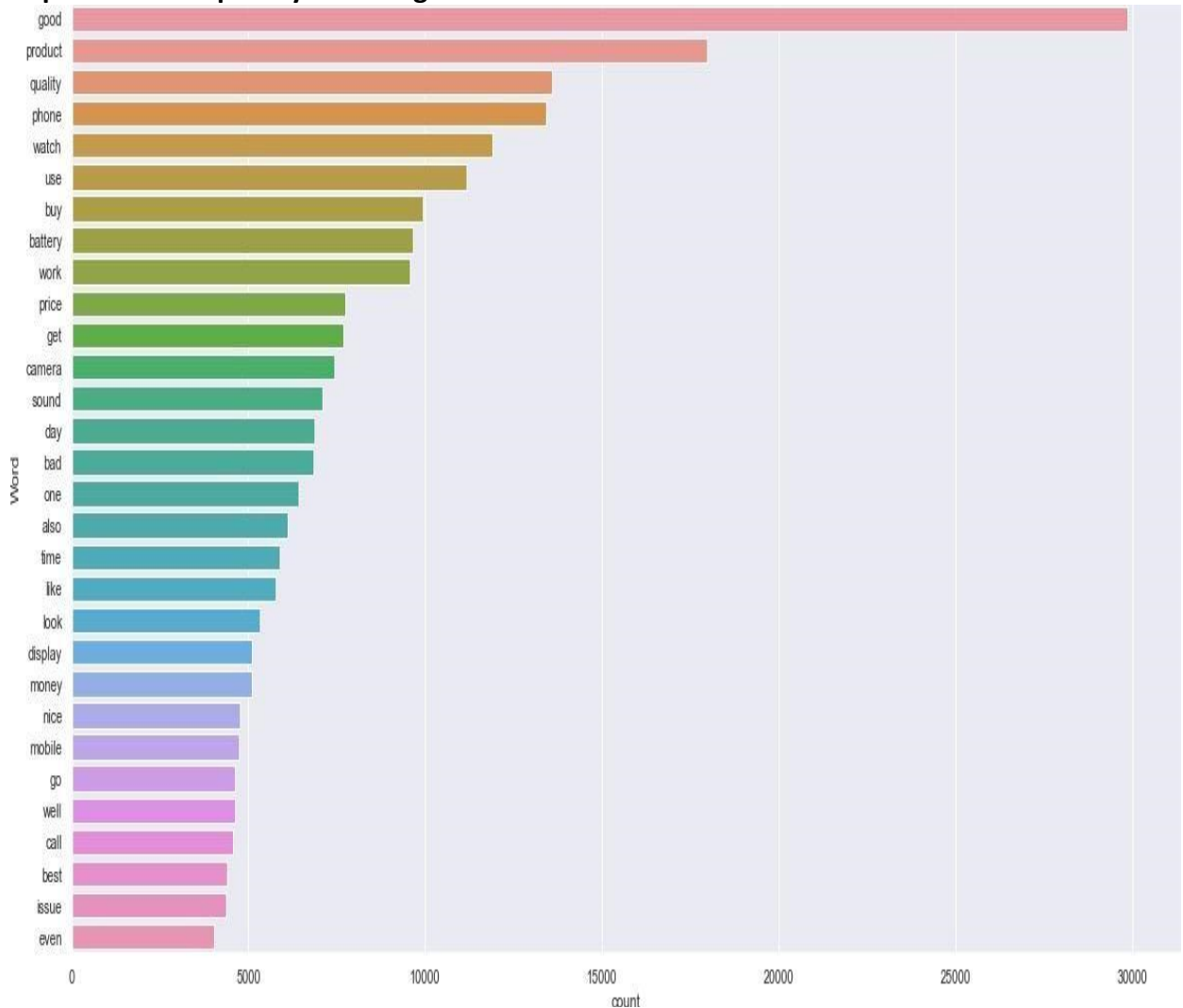


The above figure is representing count-plot for our target variable that is “Ratings”. Looking at this plot we can say that there are more numbers of reviews rated as 5 stars than others. And the reviews which are rated as 2 stars are very less in numbers when compared to others. This will cause the problem of imbalance for our model. So I have decided to select equal number of reviews from each class. I have observed that there are 7356 numbers of reviews rated as 2 stars which are least among all. So I am selecting 7356 numbers of reviews from each class as input for our model to eliminate the problem of imbalance from our data set.

```
[51] 1 #selct data from every category
2 df1 = df[df['Ratings']==1][0:7356]
3 df2 = df[df['Ratings']==2][0:7356]
4 df3 = df[df['Ratings']==3][0:7356]
5 df4 = df[df['Ratings']==4][0:7356]
6 df5 = df[df['Ratings']==5][0:7356]

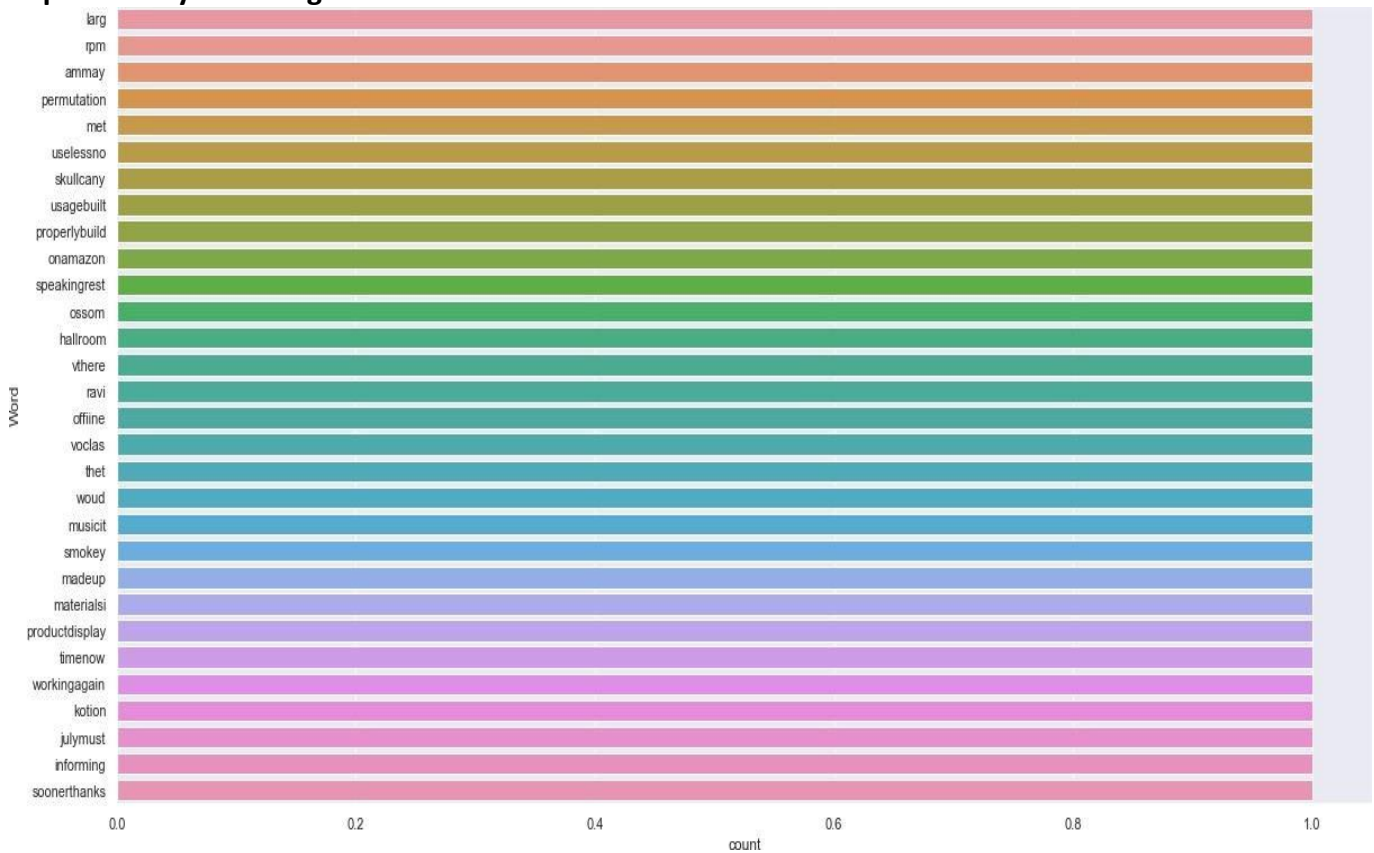
[52] 1 #Combining all the dataframes into one and shuffling them
2 df = pd.concat([df1,df2,df3,df4,df5],ignore_index=True)
3 df = df.sample(frac=1).reset_index(drop=True)
4 df
```

### Top 30 most frequently occurring words:



The above bar plot is showing top 30 most frequently occurring words in our reviews. We can see the words like 'good', 'product', 'quality' etc. are occurring more frequently.

### Top 30 Rarely occurring words:



Above figure is representing bar plot for top 30 rarely occurring words. Many of which are spelled incorrectly that's why these are occurring only once. Now using word cloud I have visualized the frequently occurring words with respect to particular rating.

**Words for rating = 1:**



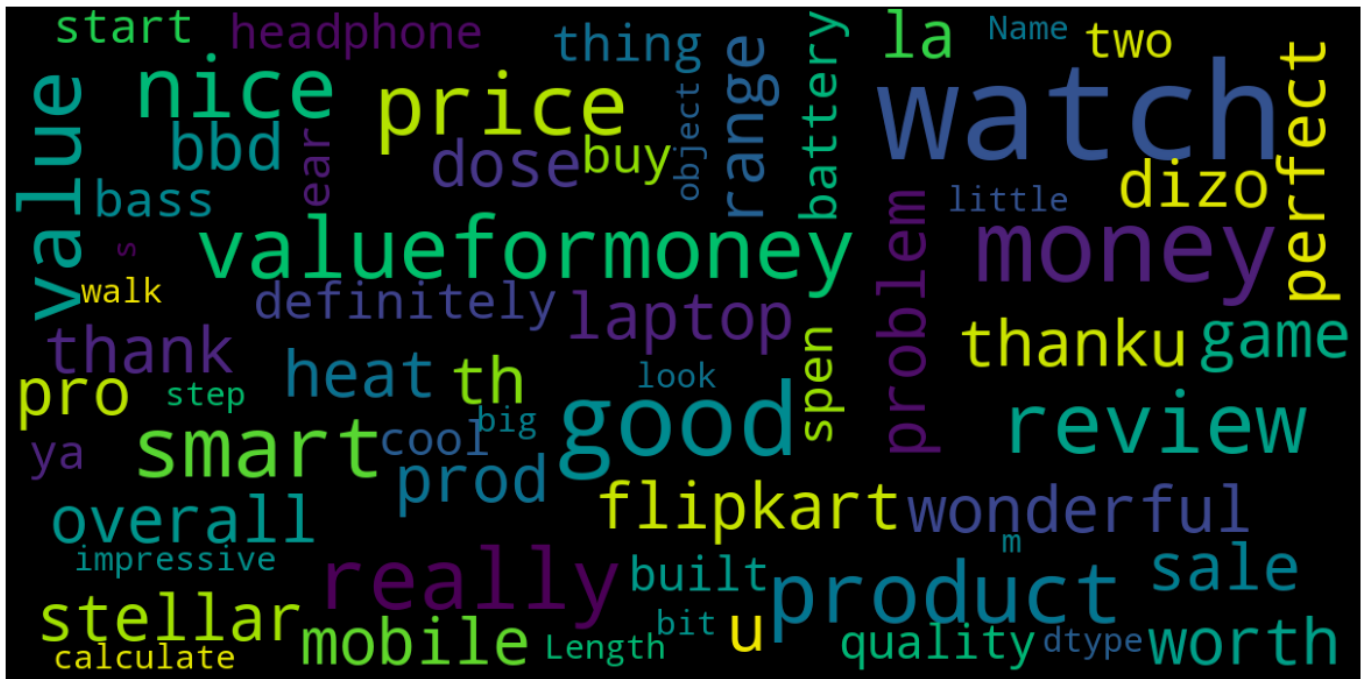
**Words for rating = 2:**



**Words for rating = 3:**



**Words for rating = 4:**



Words for rating = 5:



## Model Development and Evaluation:

As for this project we are going to predict the ratings based on the reviews given by customers this will be a classification task. For this purpose I have collected data from amazon and flipkart.

Going through various NLP steps and analyzing the data using different EDA steps I have build several models using the **Tfidf vectorizer**. Among all the different algorithms i have used Linear SVC is giving highest accuracy. Other algorithms like LGBMClassifier, XGBClassifier and RandomForestClassifier are also giving good accuracies. Considering all f1\_scores, recall and precision for different classes and cross validation score I can say the LinearSVC is giving better performance than others. So I am selecting it as best suitable algorithm for our final model. I have used following algorithms and evaluated them:

- RandomForestClassifier
  - LinearSVC
  - LogisticRegression
  - MultinomialNB
  - XGBClassifier
  - BernoulliNB
  - LightGBMClassifier
  - SGDClassifier
- 
- From all of these above models LinearSVC was giving me good performance.

## Key Metrics for success in solving problems under consideration

I have used the following metrics for evaluation:

- As this is classification problem I am using accuracy score here.
- In this case I have checked for the confusion matrix which will give clear idea about true and false predictions.
- I have checked for classification report which gives overall performance metric of any algorithm with all f1\_scores, precisions and recall scores.
- And Cross-validation score for checking the model performance for different folds.

### Hyperparameter Tuning

I have did hyperparameter tuning for LinearSVC for the parameters like 'penalty', 'loss', 'multi\_class', 'intercept\_scaling', 'dual'.

```
> 1 GCV.best_params_ #printing the best parameters found by GridSearchCV
[85]
... {'dual': True,
     'intercept_scaling': 3,
     'loss': 'hinge',
     'multi_class': 'ovr',
     'penalty': 'l2'}
```

And after doing hyper-parameter tuning I got above parameters as best suitable parameters for our final model. I have tested my final model using these parameters and got better results compared to earlier results for my final model.



## Final Model:

### Final Model

```
1 #training and testing our final model with above parameters
2 model = LinearSVC(dual = True, intercept_scaling = 2, loss = 'hinge', multi_class = 'ovr', penalty = 'l2')
3 model.fit(x_train,y_train) #fitting data to model
4 pred = model.predict(x_test)
5 accuracy = accuracy_score(y_test,pred)*100
6
7 #printing accuracy score
8 print("Accuracy Score :", accuracy)
9
10 #printing Confusion matrix
11 print(f"\nConfusion Matrix : \n {confusion_matrix(y_test,pred)}\n")
12
13 #printing Classification report
14 print(f"\nCLASSIFICATION REPORT : \n {classification_report(y_test,pred)}")
```

[87]

... Accuracy Score : 71.22349102773246

Confusion Matrix :

```
[[1416 259 118 31 18]
 [ 317 1152 233 86 22]
 [ 142 290 1137 250 61]
 [ 50 84 186 1323 193]
 [ 30 31 40 205 1521]]
```

CLASSIFICATION REPORT :

	precision	recall	f1-score	support
1	0.72	0.77	0.75	1842
2	0.63	0.64	0.64	1810
3	0.66	0.60	0.63	1880
4	0.70	0.72	0.71	1836
5	0.84	0.83	0.84	1827
accuracy			0.71	9195
macro avg	0.71	0.71	0.71	9195
weighted avg	0.71	0.71	0.71	9195

Great after doing hyperparameter tuning we have got an improved accuracy score for our final model.

## **Conclusion:**

### **Key findings of the study**

In this project I have collected data of reviews and ratings for different products from amazon.in and flipkart.com. Then I have done different text processing for reviews column and chose equal number of text from each rating class to eliminate problem of imbalance. By doing different EDA steps I have analyzed the text. We have checked frequently occurring words in our data as well as rarely occurring words. After all these steps I have built function to train and test different algorithms and using various evaluation metrics I have selected LinearSVC for our final model.

Finally by doing hyperparameter tuning we got optimum parameters for our final model. And finally we got improved accuracy score for our final model.

### **Limitations of this work and scope for the future work**

As we know the content of text in reviews is totally depends on the reviewer and they may rate differently which is totally depends on that particular person. So it is difficult to predict ratings based on the reviews with higher accuracies.

Still we can improve our accuracy by fetching more data and by doing extensive hyperparameter tuning.

