# NAME OF THE PROJECT

# Product Review Rating Using NLP

**Submitted by:**
Arshi Maitra

# Acknowledgement:

This project has been completed with the help of training documents and live classes recordings from Data Trained Education. Few help on coding have also been taken from few data science websites like Toward Data Science, Geek for Geeks, Stack Overflow. All data related to product rating has been taken from FlipKart under Headphones, Camera, Speakers, Mobiles, Iphone and TV. I would also like to extend my gratitude to FlipRobo team who have given me this wonderful opportunity opportunity to work on this amazing dataset.

# INTRODUCTION

Machine Learning (ML) is a field of Artificial Intelligence where data-driven algorithms learn patterns by getting exposed to relevant data. ML has gained massive importance in the field of Natural Language Processing (NLP), i.e. interpreting human language. In this article, we will focus on the use of both ML and NLP in predicting the ratings of user reviews.

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.

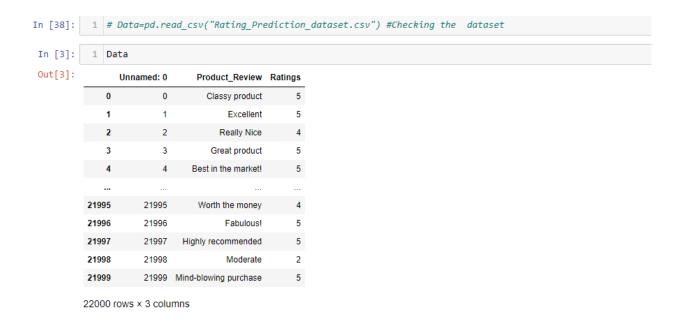Therefore the main goal of this project is to determine the following:

• To predict the rating of the product by looking at their review

**The following data has been taken from flipkart.com. Here the ratings of various electronic products have been scraped around 22,000. We will be trying to analyze the pattern of words as well as the ratings prediction from the scraped dataset.**

# Analytical Problem Framing

# Web Scraping

Before importing the necessary libraries for data analysis and prediction, the data has been first scrapped from the website www.flipkart.com by using the Selenium method from Web Scraping. The data is saved in CSV which is later on used for analysis:

```
In [38]:   1   # Data=pd.read_csv("Rating_Prediction_dataset.csv") #Checking the  dataset
```

```
In [3]:   1   Data
```

Out[3]:

|       | Unnamed: 0 | Product_Review | Ratings |
|-------|------------|----------------|---------|
| 0     | 0          | Classy product | 5       |
| 1     | 1          | Excellent      | 5       |
| 2     | 2          | Really Nice    | 4       |
| 3     | 3          | Great product  | 5       |
| 4     | 4          | Best in the market! | 5  |
| ...   | ...        | ...            | ...     |
| 21995 | 21995      | Worth the money | 4      |
| 21996 | 21996      | Fabulous!      | 5       |
| 21997 | 21997      | Highly recommended | 5   |
| 21998 | 21998      | Moderate       | 2       |
| 21999 | 21999      | Mind-blowing purchase | 5 |

22000 rows × 3 columns

# Importing Necessary Libraries
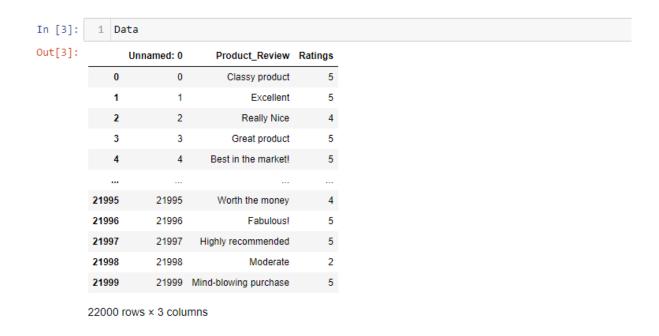
## ML:

```
In [37]:   1  import numpy as np
           2  import pandas as pd
           3  import seaborn as sns
           4  import matplotlib.pyplot as plt
           5  from sklearn.preprocessing import StandardScaler
           6  from statsmodels.stats.outliers_influence import variance_inflation_factor
           7  from sklearn.model_selection import cross_val_score
           8  from sklearn.model_selection import train_test_split
           9  from sklearn.linear_model import Ridge
          10  from sklearn.linear_model import Lasso
          11  from sklearn.metrics import r2_score
          12  from sklearn.tree import DecisionTreeRegressor
          13  from sklearn.ensemble import RandomForestRegressor
          14  from sklearn.preprocessing import LabelEncoder
          15  from sklearn.model_selection import GridSearchCV
          16  from sklearn.linear_model import LogisticRegression
          17  from sklearn.ensemble import RandomForestClassifier
          18  from sklearn.neighbors import KNeighborsClassifier
          19  from sklearn.ensemble import VotingClassifier
          20  from sklearn.ensemble import GradientBoostingClassifier
          21  from sklearn.tree import DecisionTreeClassifier
          22  from sklearn.svm import SVC
          23  from sklearn.metrics import confusion_matrix, classification_report
          24  from sklearn.metrics import accuracy_score,confusion_matrix,roc_curve, roc_auc_score
          25  from prettytable import PrettyTable
          26  import warnings
          27  warnings.filterwarnings('ignore') #Importing the necessary libraries
```

## NLP:

```
1  #Importing required libraries
2  import re
3  import string
4  import nltk
5  from nltk.corpus import stopwords
6  from nltk.tokenize import word_tokenize
7  from nltk.stem import SnowballStemmer, WordNetLemmatizer
8  from sklearn.feature_extraction.text import CountVectorizer, TfidfVectorizer
9  from wordcloud import WordCloud
```

# Data preprocessing/Data Cleaning:

## Loading the dataset

```
In [3]:    1  Data
```

Out[3]:

|  | Unnamed: 0 | Product_Review | Ratings |
|---|---|---|---|
| 0 | 0 | Classy product | 5 |
| 1 | 1 | Excellent | 5 |
| 2 | 2 | Really Nice | 4 |
| 3 | 3 | Great product | 5 |
| 4 | 4 | Best in the market! | 5 |
| ... | ... | ... | ... |
| 21995 | 21995 | Worth the money | 4 |
| 21996 | 21996 | Fabulous! | 5 |
| 21997 | 21997 | Highly recommended | 5 |
| 21998 | 21998 | Moderate | 2 |
| 21999 | 21999 | Mind-blowing purchase | 5 |

22000 rows × 3 columns

## Checking the number of rows and columns

```
In [7]:    1  print('No. of Rows :',Data.shape[0])
           2  print('No. of Columns :',Data.shape[1])
```

```
No. of Rows : 22000
No. of Columns : 2
```

**Checking the null values from the  data:**

```
In [8]:   1  Data.isnull().sum() # Checking the null values

Out[8]:  Product_Review      0
         Ratings             0
         dtype: int64
```

**Dropping the unnecessary column:**

```
1  # Dropping unnecssary index column Unnamed:0
2  Data.drop('Unnamed: 0',axis=1,inplace=True)
```

# Visualization:

**Checking Rating Distribution:**

```
1  sns.countplot(Data['Ratings']) # Checking the distribution of ratings
```
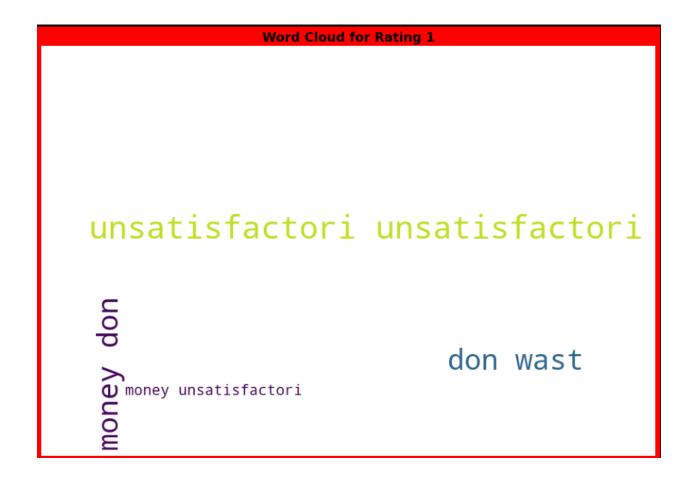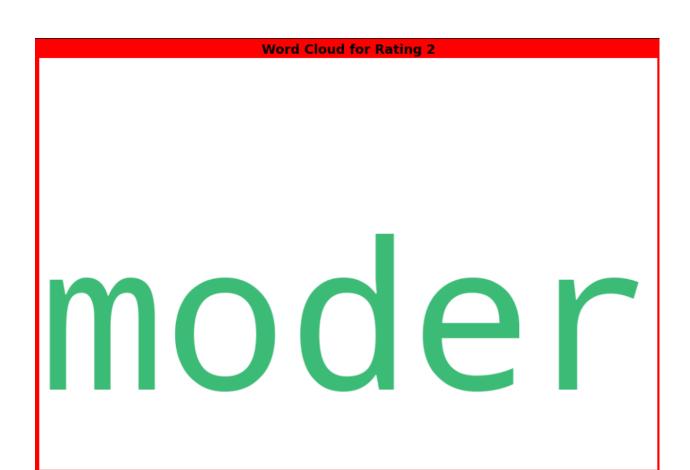
```
<AxesSubplot:xlabel='Ratings', ylabel='count'>
```
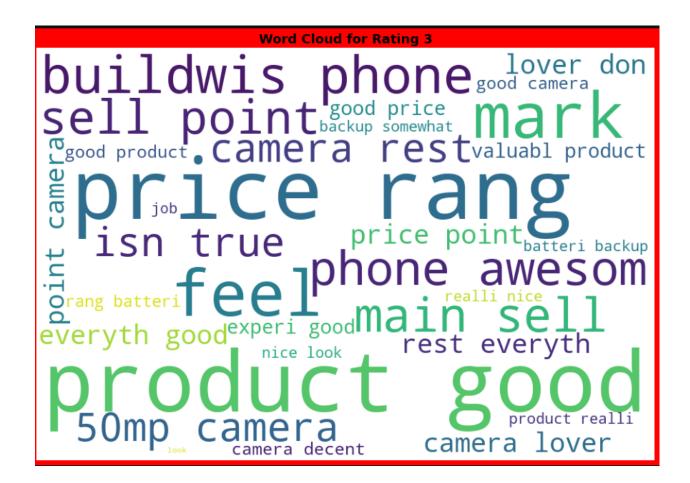
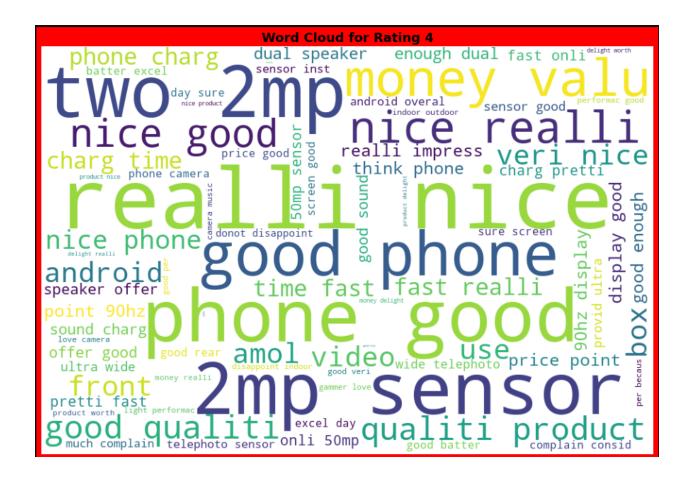**Inference: The highest number of rating is 5 while rating 2 is the lowest.**

**WordCloud:**

Word clouds or tag clouds are graphical representations of word frequency that give greater prominence to words that appear more frequently in a source text. The larger the word in the visual the more common the word was in the document

**Word Cloud for Rating 1**

unsatisfactori unsatisfactori

money don

money unsatisfactori

don wast

**Word Cloud for Rating 2**

moder

Word Cloud for Rating 3

Word Cloud for Rating 4

# Word Cloud for Rating 5

# Model/s Development and Evaluation

# (Logistic Regression, Decision Tree Classifier, Random Forest Classifier)

# Logistic Regression:

## Training and Testing the model

```
In [39]:    1  LR=LogisticRegression()
```

```
In [41]:    1  X_train,X_test,y_train,y_test=train_test_split(X,Y,test_size=0.20,random_state=58)
            2  LR.fit(X_train,y_train)
            3  pred_test=LR.predict(X_test) # Testing the prediction of test data
```

```
In [42]:    1  LR_accuracy=accuracy_score(y_test,pred_test)*100
            2  LR_accuracy
```

```
Out[42]:  100.0
```

## Confusion Matrix:

```
array([[  86,     0,     0,     0,     0],
       [   0,    50,     0,     0,     0],
       [   0,     0,   228,     0,     0],
       [   0,     0,     0,   963,     0],
       [   0,     0,     0,     0, 3073]], dtype=int64)
```

# Decision Tree Classifier

## Training and Testing the model

```
1  DT=DecisionTreeClassifier()
```

```
1  X_train,X_test,y_train,y_test=train_test_split(X,Y,test_size=0.20,random_state=58)
2  DT.fit(X_train,y_train)
3  pred_test=DT.predict(X_test) # Testing the prediction of test data
```

```
1  DT_accuracy=accuracy_score(y_test,pred_test)*100
2  DT_accuracy
```

```
100.0
```

**Confusion Matrix:**

```
array([[  86,    0,    0,    0,    0],
       [   0,   50,    0,    0,    0],
       [   0,    0,  228,    0,    0],
       [   0,    0,    0,  963,    0],
       [   0,    0,    0,    0, 3073]], dtype=int64)
```

# Random Forest Classifier

## Training and Testing the model

```
1  RF=RandomForestClassifier()
```

```
1  X_train,X_test,y_train,y_test=train_test_split(X,Y,test_size=0.20,random_state=58)
2  RF.fit(X_train,y_train)
3  pred_test=RF.predict(X_test) # Testing the prediction of test data
```

```
1  RF_accuracy=accuracy_score(y_test,pred_test)*100
2  RF_accuracy
```

100.0

## Confusion Matrix:

```
array([[  86,    0,    0,    0,    0],
       [   0,   50,    0,    0,    0],
       [   0,    0,  228,    0,    0],
       [   0,    0,    0,  963,    0],
       [   0,    0,    0,    0, 3073]], dtype=int64)
```

## Hyperparameter Tuning:

As all the scores came down to 100%, hence the hypertuning parameter is not performed. Random Forest Classifier is chosen since It is based on the concept of ensemble learning, which is a process of combining multiple classifiers to solve a complex problem and to improve the performance of the model.

# Conclusion

## Learning Outcomes:

1. Hands on chance to enhance  web scraping skills

2. In this project, I was able to learn various NLP techniques like Stemming, Lemmatization and removal of stop words.

3. This project has stressed the importance of sampling, modeling and predicting data.

## Limitations:

1. More input features could be scraped.
2. There is scope for application of advance deep learning NLP tool.

# Thank You!