

Zomato Dataset Exploratory Data Analysis

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import warnings
warnings.filterwarnings('ignore')
```

```
In [2]: # "D:\zomato.csv"
df=pd.read_csv("D://zomato.csv",encoding='latin_1')
df.head()
```

Out[2]:

	Restaurant ID	Restaurant Name	Country Code	City	Address	Locality	Locality Verbose	Longitude	Latitude	Cuisines	...	Currency	Has Table booking	Has Online delivery	deliv
0	6317637	Le Petit Souffle	162	Makati City	Third Floor, Century City Mall, Kalayaan Avenu...	Century City Mall, Poblacion, Makati City	Century City Mall, Poblacion, Makati City, Mak...	121.027535	14.565443	French, Japanese, Desserts	...	Botswana Pula(P)	Yes	No	
1	6304287	Izakaya Kikufuji	162	Makati City	Little Tokyo, 2277 Chino Roces Avenue, Legaspi...	Little Tokyo, Legaspi Village, Makati City	Little Tokyo, Legaspi Village, Makati City, Ma...	121.014101	14.553708	Japanese	...	Botswana Pula(P)	Yes	No	
2	6300002	Heat - Edsa Shangri-La	162	Mandaluyong City	Edsa Shangri-La, 1 Garden Way, Ortigas, Mandal...	Edsa Shangri-La, Ortigas, Mandaluyong City	Edsa Shangri-La, Ortigas, Mandaluyong City, Ma...	121.056831	14.581404	Seafood, Asian, Filipino, Indian	...	Botswana Pula(P)	Yes	No	
3	6318506	Ooma	162	Mandaluyong City	Third Floor, Mega Fashion Hall, SM Megamall, O...	SM Megamall, Ortigas, Mandaluyong City	SM Megamall, Ortigas, Mandaluyong City, Mandal...	121.056475	14.585318	Japanese, Sushi	...	Botswana Pula(P)	No	No	
4	6314302	Sambo Kojin	162	Mandaluyong City	Third Floor, Mega Atrium, SM Megamall, Ortigas...	SM Megamall, Ortigas, Mandaluyong City	SM Megamall, Ortigas, Mandaluyong City, Mandal...	121.057508	14.584450	Japanese, Korean	...	Botswana Pula(P)	Yes	No	

5 rows × 21 columns

```
In [3]: df.columns
```

```
Out[3]: Index(['Restaurant ID', 'Restaurant Name', 'Country Code', 'City', 'Address',
'Locality', 'Locality Verbose', 'Longitude', 'Latitude', 'Cuisines',
'Average Cost for two', 'Currency', 'Has Table booking',
'Has Online delivery', 'Is delivering now', 'Switch to order menu',
'Price range', 'Aggregate rating', 'Rating color', 'Rating text',
'Votes'],
dtype='object')
```

```
In [4]: df.shape
```

```
Out[4]: (9551, 21)
```

```
In [5]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 9551 entries, 0 to 9550
Data columns (total 21 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Restaurant ID          9551 non-null   int64
1   Restaurant Name        9551 non-null   object
2   Country Code           9551 non-null   int64
3   City                   9551 non-null   object
4   Address                9551 non-null   object
5   Locality               9551 non-null   object
6   Locality Verbose       9551 non-null   object
7   Longitude              9551 non-null   float64
8   Latitude               9551 non-null   float64
9   Cuisines               9542 non-null   object
10  Average Cost for two   9551 non-null   int64
11  Currency               9551 non-null   object
12  Has Table booking      9551 non-null   object
13  Has Online delivery    9551 non-null   object
14  Is delivering now      9551 non-null   object
15  Switch to order menu   9551 non-null   object
16  Price range            9551 non-null   int64
17  Aggregate rating       9551 non-null   float64
18  Rating color           9551 non-null   object
19  Rating text            9551 non-null   object
20  Votes                  9551 non-null   int64
dtypes: float64(3), int64(5), object(13)
memory usage: 1.5+ MB
```

```
In [6]: df.describe()
```

Out[6]:

	Restaurant ID	Country Code	Longitude	Latitude	Average Cost for two	Price range	Aggregate rating	Votes
count	9.551000e+03	9551.000000	9551.000000	9551.000000	9551.000000	9551.000000	9551.000000	9551.000000
mean	9.051128e+06	18.365616	64.126574	25.854381	1199.210763	1.804837	2.666370	156.909748
std	8.791521e+06	56.750546	41.467058	11.007935	16121.183073	0.905609	1.516378	430.169145
min	5.300000e+01	1.000000	-157.948486	-41.330428	0.000000	1.000000	0.000000	0.000000
25%	3.019625e+05	1.000000	77.081343	28.478713	250.000000	1.000000	2.500000	5.000000
50%	6.004089e+06	1.000000	77.191964	28.570469	400.000000	2.000000	3.200000	31.000000
75%	1.835229e+07	1.000000	77.282006	28.642758	700.000000	2.000000	3.700000	131.000000
max	1.850065e+07	216.000000	174.832089	55.976980	800000.000000	4.000000	4.900000	10934.000000

In Data Analysis What All Things We Do

- 1. Missing Values
- 2. Explore About the Numerical Variables
- 3. Explore About categorical variables
- 4. Finding Relationship between features

```
In [7]: df.isna().sum()
```

Out[7]:

```
Restaurant ID      0
Restaurant Name     0
Country Code       0
City               0
Address            0
Locality           0
Locality Verbose   0
Longitude          0
Latitude           0
Cuisines           9
Average Cost for two 0
Currency           0
Has Table booking   0
Has Online delivery 0
Is delivering now   0
Switch to order menu 0
Price range        0
Aggregate rating    0
Rating color       0
Rating text        0
Votes              0
dtype: int64
```

```
In [8]: [features for features in df.columns if df[features].isnull().sum()>0]
```

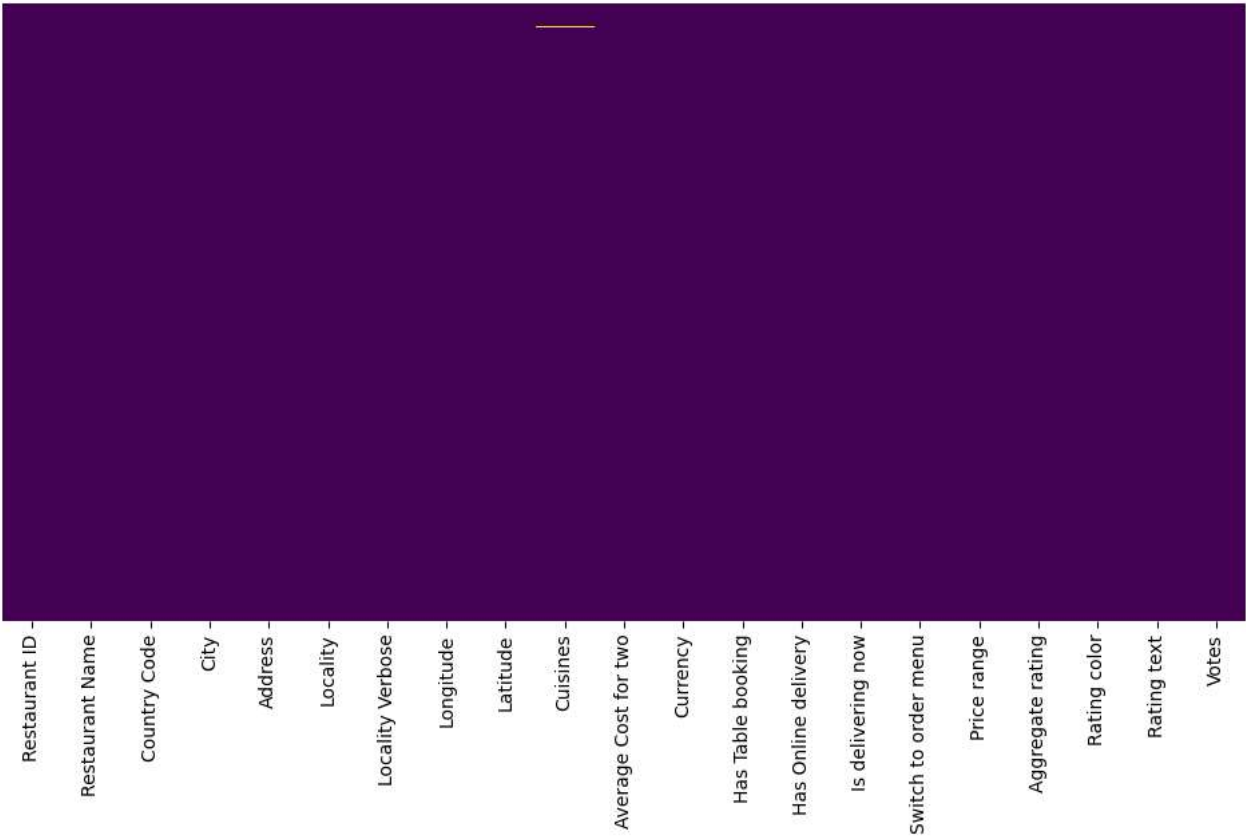
```
Out[8]: ['Cuisines']
```

```
In [9]: [features for features in df.columns if df[features].isnull().sum()>0]
```

```
Out[9]: ['Cuisines']
```

```
In [10]: import matplotlib
matplotlib.rcParams['figure.figsize']=(12,6)
sns.heatmap(df.isnull(),yticklabels=False,cbar=False,cmap='viridis')
```

```
Out[10]: <AxesSubplot:>
```



```
In [11]: df_country=pd.read_excel("D://Country-Code.xlsx")
```

```
In [12]: df_country
```

```
Out[12]:
```

	Country Code	Country
0	1	India
1	14	Australia
2	30	Brazil
3	37	Canada
4	94	Indonesia
5	148	New Zealand
6	162	Phillipines
7	166	Qatar
8	184	Singapore
9	189	South Africa
10	191	Sri Lanka
11	208	Turkey
12	214	UAE
13	215	United Kingdom
14	216	United States

In [13]:

df.columns

Out[13]:

Index(['Restaurant ID', 'Restaurant Name', 'Country Code', 'City', 'Address',
 'Locality', 'Locality Verbose', 'Longitude', 'Latitude', 'Cuisines',
 'Average Cost for two', 'Currency', 'Has Table booking',
 'Has Online delivery', 'Is delivering now', 'Switch to order menu',
 'Price range', 'Aggregate rating', 'Rating color', 'Rating text',
 'Votes'],
 dtype='object')

In [14]:

final_df=pd.merge(df,df_country,on='Country Code',how='left')

In [15]:

final_df.head(2)

Out[15]:

	Restaurant ID	Restaurant Name	Country Code	City	Address	Locality	Locality Verbose	Longitude	Latitude	Cuisines	...	Has Table booking	Has Online delivery	Is delivering now	Switch to order menu	Price range
0	6317637	Le Petit Souffle	162	Makati City	Third Floor, Century City Mall, Kalayaan Avenu...	Century City Mall, Poblacion, Makati City	Century City Mall, Poblacion, Makati City, Mak...	121.027535	14.565443	French, Japanese, Desserts	...	Yes	No	No	No	3
1	6304287	Izakaya Kikufuji	162	Makati City	Little Tokyo, 2277 Chino Roces Avenue, Legaspi...	Little Tokyo, Legaspi Village, Makati City	Little Tokyo, Legaspi Village, Makati City, Ma...	121.014101	14.553708	Japanese	...	Yes	No	No	No	3

2 rows × 22 columns

In [16]:

To check Data Types

final_df.dtypes

Out[16]:

Restaurant ID int64
Restaurant Name object
Country Code int64
City object
Address object
Locality object
Locality Verbose object
Longitude float64
Latitude float64
Cuisines object
Average Cost for two int64
Currency object
Has Table booking object
Has Online delivery object
Is delivering now object
Switch to order menu object
Price range int64
Aggregate rating float64
Rating color object
Rating text object
Votes int64
Country object
dtype: object

In [17]:

final_df.Country.value_counts()

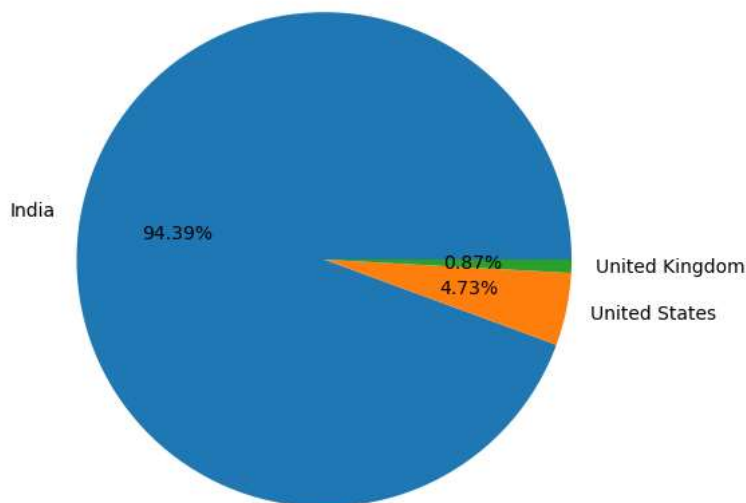
Out[17]:

India 8652
United States 434
United Kingdom 80
Brazil 60
UAE 60
South Africa 60
New Zealand 40
Turkey 34
Australia 24
Phillipines 22
Indonesia 21
Singapore 20
Qatar 20
Sri Lanka 20
Canada 4
Name: Country, dtype: int64

```
In [18]: country_names=final_df.Country.value_counts().index
```

```
In [19]: country_val=final_df.Country.value_counts().values
```

```
In [39]: ## Pie Chart-- Top 3 countries that uses zomato
plt.pie(country_val[:3],labels=country_names[:3],autopct='%1.2f%%')
plt.show()
```



Observation: Most of the orders from the India after that USA and then UK.

```
In [21]: final_df.columns
```

```
Out[21]: Index(['Restaurant ID', 'Restaurant Name', 'Country Code', 'City', 'Address',
              'Locality', 'Locality Verbose', 'Longitude', 'Latitude', 'Cuisines',
              'Average Cost for two', 'Currency', 'Has Table booking',
              'Has Online delivery', 'Is delivering now', 'Switch to order menu',
              'Price range', 'Aggregate rating', 'Rating color', 'Rating text',
              'Votes', 'Country'],
              dtype='object')
```

```
In [22]: ratings=final_df.groupby(['Aggregate rating', 'Rating color', 'Rating text']).size().reset_index().rename(columns={0: 'Rating Count'}
```

In [23]: ratings

Out[23]:

	Aggregate rating	Rating color	Rating text	Rating Count
0	0.0	White	Not rated	2148
1	1.8	Red	Poor	1
2	1.9	Red	Poor	2
3	2.0	Red	Poor	7
4	2.1	Red	Poor	15
5	2.2	Red	Poor	27
6	2.3	Red	Poor	47
7	2.4	Red	Poor	87
8	2.5	Orange	Average	110
9	2.6	Orange	Average	191
10	2.7	Orange	Average	250
11	2.8	Orange	Average	315
12	2.9	Orange	Average	381
13	3.0	Orange	Average	468
14	3.1	Orange	Average	519
15	3.2	Orange	Average	522
16	3.3	Orange	Average	483
17	3.4	Orange	Average	498
18	3.5	Yellow	Good	480
19	3.6	Yellow	Good	458
20	3.7	Yellow	Good	427
21	3.8	Yellow	Good	400
22	3.9	Yellow	Good	335
23	4.0	Green	Very Good	266
24	4.1	Green	Very Good	274
25	4.2	Green	Very Good	221
26	4.3	Green	Very Good	174
27	4.4	Green	Very Good	144
28	4.5	Dark Green	Excellent	95
29	4.6	Dark Green	Excellent	78
30	4.7	Dark Green	Excellent	42
31	4.8	Dark Green	Excellent	25
32	4.9	Dark Green	Excellent	61

Observation

1. when the rating in between 4.5 to 4.9----> Excellent
2. when the rating in between 4.0 to 4.4----> Very Good
3. when the rating in between 3.5 to 3.9----> Good
4. when the rating in between 2.5 to 3.4----> Average
5. when the rating in between 1.8 to 2.4----> Poor

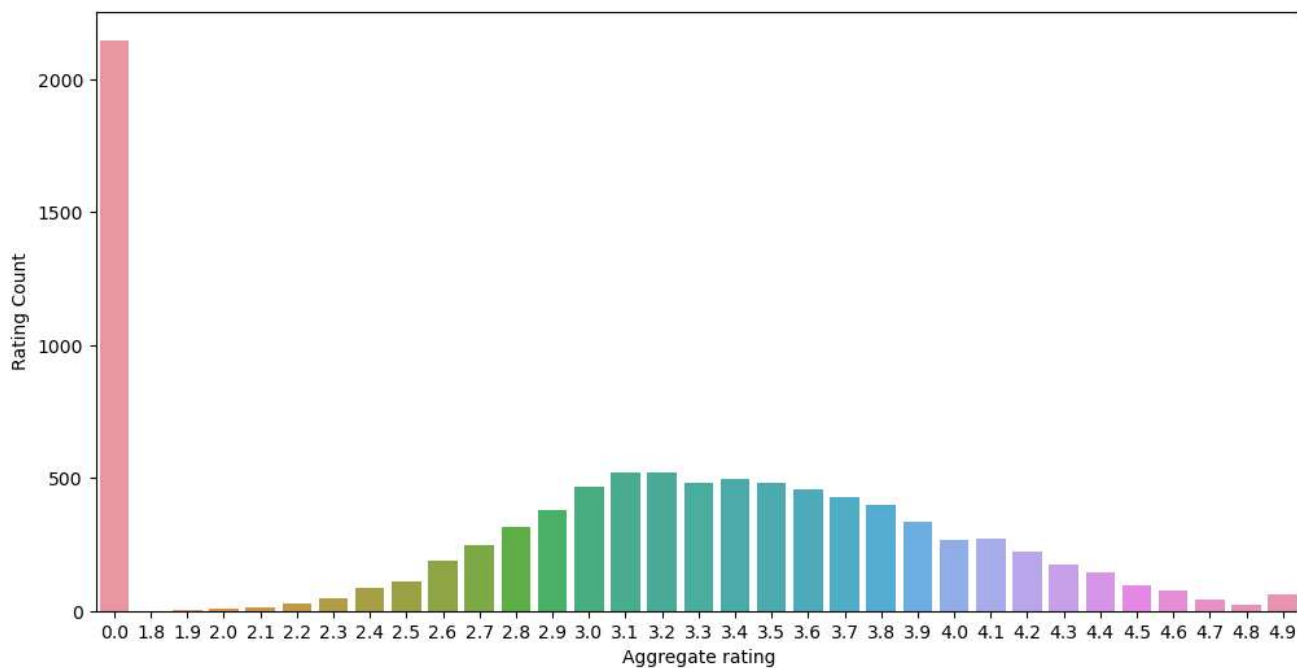
In [24]: ratings.head()

Out[24]:

	Aggregate rating	Rating color	Rating text	Rating Count
0	0.0	White	Not rated	2148
1	1.8	Red	Poor	1
2	1.9	Red	Poor	2
3	2.0	Red	Poor	7
4	2.1	Red	Poor	15

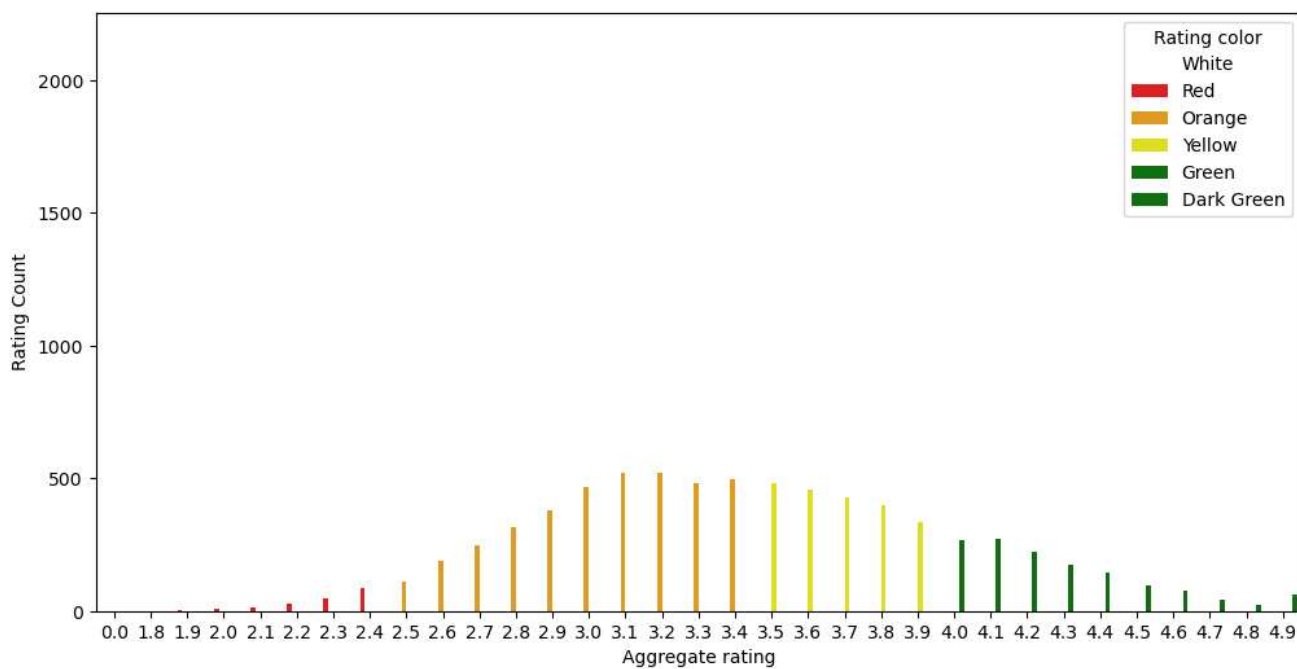
```
In [25]: plt.figure(figsize=(12,6))
sns.barplot(x='Aggregate rating',y='Rating Count',data=ratings)
```

```
Out[25]: <AxesSubplot:xlabel='Aggregate rating', ylabel='Rating Count'>
```



```
In [26]: sns.barplot(x='Aggregate rating',y='Rating Count',hue='Rating color',data=ratings,palette=['white','Red','Orange','Yellow','Green','Dark Green'])
```

```
Out[26]: <AxesSubplot:xlabel='Aggregate rating', ylabel='Rating Count'>
```

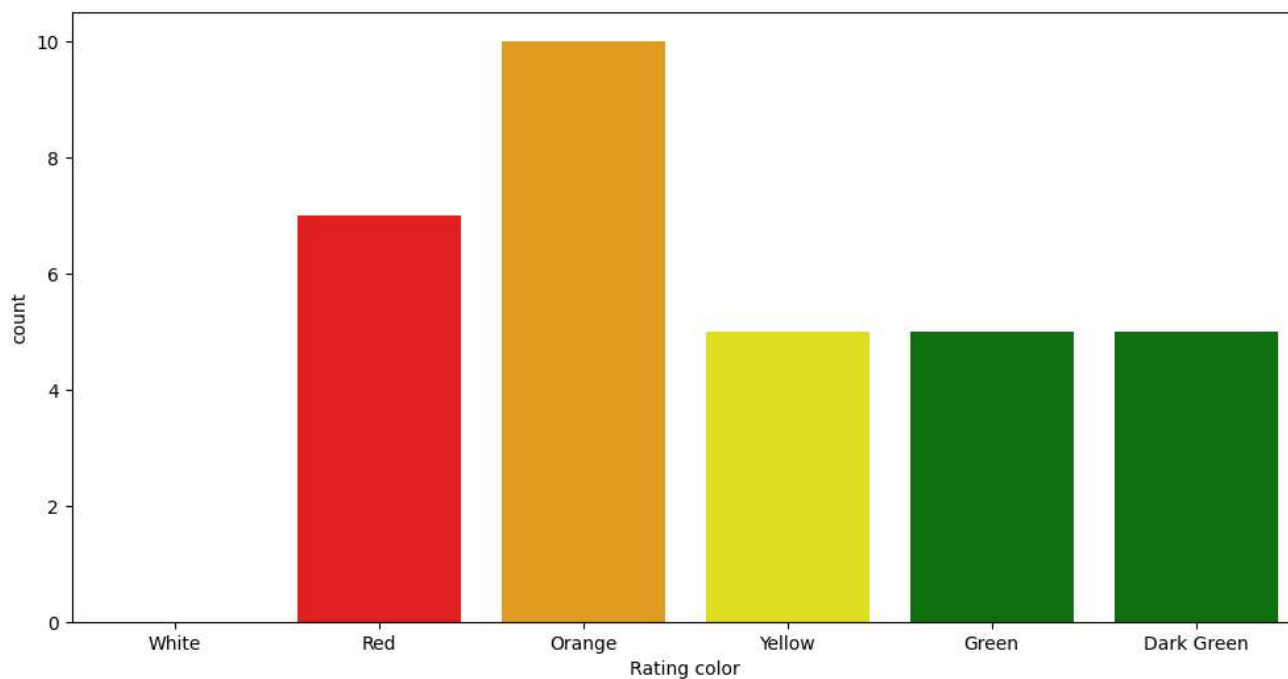


Observation:

1. Not Rated count is very high
2. Maximum number of rating are between 2.5 to 3.4

```
In [27]: ## Count plot
sns.countplot(x='Rating color',data=ratings,palette=['White','Red','Orange','Yellow','Green','Green'])
```

```
Out[27]: <AxesSubplot:xlabel='Rating color', ylabel='count'>
```



```
In [28]: ### Find the countries name that has given 0 rating
final_df[final_df['Rating color']=='White'].groupby('Country').size().reset_index()
```

```
Out[28]:
```

	Country	0
0	Brazil	5
1	India	2139
2	United Kingdom	1
3	United States	3

Observations : Maximum number of ratings are from the indian customers

```
In [29]: final_df['Currency'].value_counts()
```

```
Out[29]: Indian Rupees(Rs.)      8652
Dollar($)          482
Pounds( £)         80
Brazilian Real(R$)  60
Emirati Diram(AED)  60
Rand(R)            60
NewZealand($)       40
Turkish Lira(TL)    34
Botswana Pula(P)    22
Indonesian Rupiah(IDR) 21
Qatari Rial(QR)     20
Sri Lankan Rupee(LKR) 20
Name: Currency, dtype: int64
```



```
In [30]: ##find out which currency is used by which country?
final_df[['Country', 'Currency']].groupby(['Country', 'Currency']).size().reset_index()
```

Out[30]:

	Country	Currency	0
0	Australia	Dollar(\$)	24
1	Brazil	Brazilian Real(R\$)	60
2	Canada	Dollar(\$)	4
3	India	Indian Rupees(Rs.)	8652
4	Indonesia	Indonesian Rupiah(IDR)	21
5	New Zealand	NewZealand(\$)	40
6	Phillipines	Botswana Pula(P)	22
7	Qatar	Qatari Rial(QR)	20
8	Singapore	Dollar(\$)	20
9	South Africa	Rand(R)	60
10	Sri Lanka	Sri Lankan Rupee(LKR)	20
11	Turkey	Turkish Lira(TL)	34
12	UAE	Emirati Diram(AED)	60
13	United Kingdom	Pounds(£)	80
14	United States	Dollar(\$)	434

```
In [31]: ## Which countrys do have online deliveries option
final_df[['Country', 'Has Online delivery']].groupby(['Country', 'Has Online delivery']).size().reset_index()
```

Out[31]:

	Country	Has Online delivery	0
0	Australia	No	24
1	Brazil	No	60
2	Canada	No	4
3	India	No	6229
4	India	Yes	2423
5	Indonesia	No	21
6	New Zealand	No	40
7	Phillipines	No	22
8	Qatar	No	20
9	Singapore	No	20
10	South Africa	No	60
11	Sri Lanka	No	20
12	Turkey	No	34
13	UAE	No	32
14	UAE	Yes	28
15	United Kingdom	No	80
16	United States	No	434

Observations:

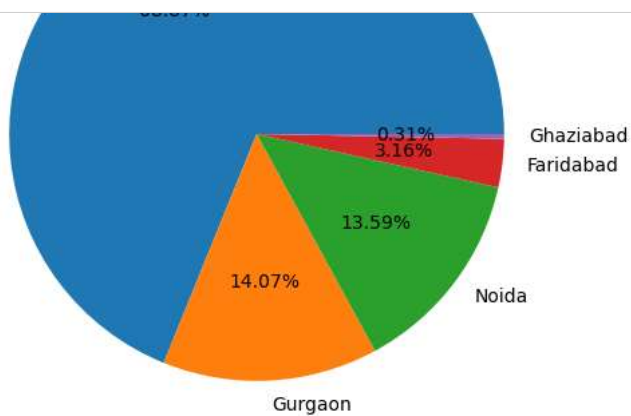
1. Online Deliveries are available in India and UAE

```
In [32]: final_df['City'].value_counts().index
```

Out[32]: Index(['New Delhi', 'Gurgaon', 'Noida', 'Faridabad', 'Ghaziabad',
'Bhubaneshwar', 'Amritsar', 'Ahmedabad', 'Lucknow', 'Guwahati',
...,
'Ojo Caliente', 'Montville', 'Monroe', 'Miller', 'Middleton Beach',
'Panchkula', 'Mc Millan', 'Mayfield', 'Macedon', 'Vineland Station'],
dtype='object', length=141)

```
In [33]: city_values=final_df.City.value_counts().values
city_labels=final_df.City.value_counts().index
```

```
In [38]: plt.pie(city_values[:5],labels=city_labels[:5],autopct='%1.2f%%')  
plt.show()
```



Assignment

Find out top 10 cuisines

```
In [35]: final_df['Cuisines'].value_counts().head(10)
```

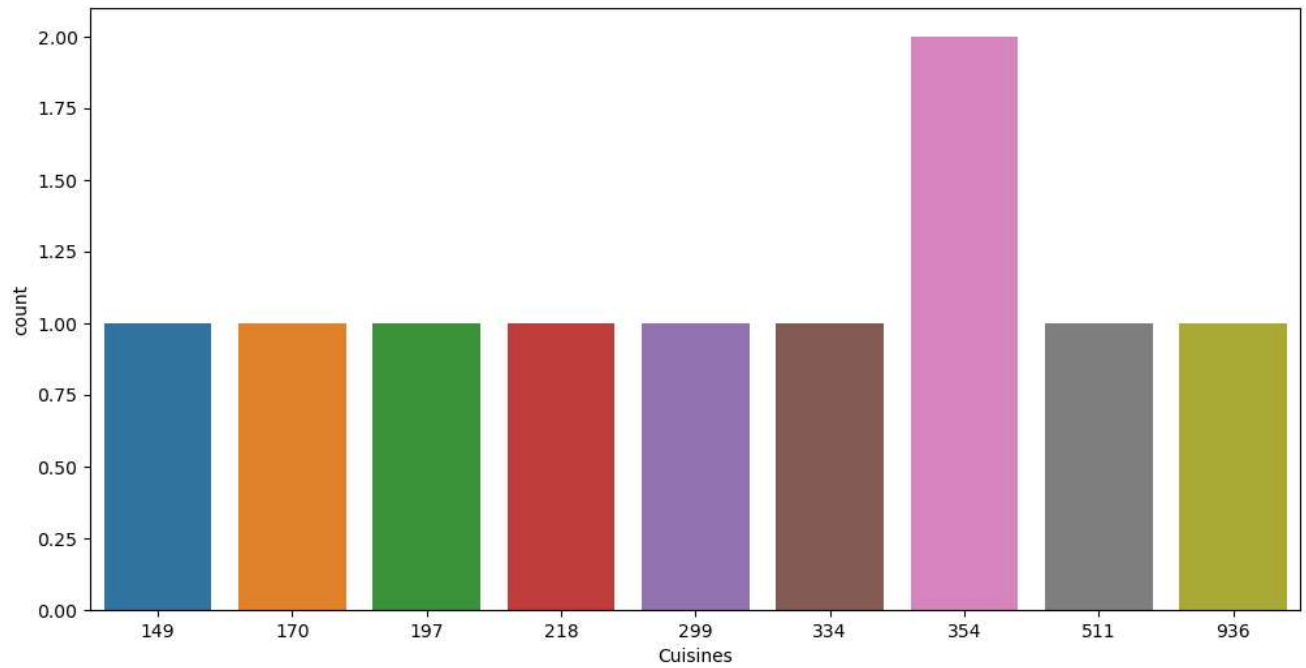
```
Out[35]: North Indian          936  
North Indian, Chinese        511  
Chinese                      354  
Fast Food                   354  
North Indian, Mughlai        334  
Cafe                        299  
Bakery                      218  
North Indian, Mughlai, Chinese 197  
Bakery, Desserts             170  
Street Food                  149  
Name: Cuisines, dtype: int64
```

```
In [36]: final_df.Cuisines.value_counts()[:5]
```

```
Out[36]: North Indian          936  
North Indian, Chinese        511  
Chinese                      354  
Fast Food                   354  
North Indian, Mughlai        334  
Name: Cuisines, dtype: int64
```

```
In [37]: sns.countplot(df['Cuisines'].value_counts()[:10])
```

```
Out[37]: <AxesSubplot:xlabel='Cuisines', ylabel='count'>
```



```
In [ ]:
```