

Analysis of Swiggy sales and customer oriented data

Intoduction Our project aims to analyze and visualize the geographical distribution of restaurants listed in a dataset across various cities. By leveraging Python programming and data visualization techniques, we seek to gain insights into the spatial distribution of restaurants and their presence in different urban centers.4ur project aims to analyze and visualize the geographical distribution of restaurants listed in a dataset across various cities. By leveraging Python programming and data visualization techniques, we seek to gain insights into the spatial distribution of restaurants and their presence in different urban centers.

About the Dataset

link: https://drive.google.com/file/d/13M8upJ0gFkaJcQY0EjLAy45c32Lb2gK/view?usp=drive_link

Dataset Analysis

Import Libraries

```
In [1]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
import geopandas as gpd
```

Import dataset

```
In [2]: df = pd.read_csv('swiggy.csv')
```

Display all column of dataset with 5 rows

```
In [3]: df.head(5)
```

	id	name	city	rating	rating_count	cost	lic_no	link	address	menu
0	567335	AB FOODS POINT	Abdohar	4.0	Too Few Ratings	200.0	2.21E+13	https://www.swiggy.com/restaurants/ab-foods-point-1	AB FOODS POINT, NEAR RISH NARANG DENTAL CLINIC, ...	Menu/567335.json
1	531342	Janta Sweet House	Abdohar	4.4	50+ ratings	200.0	1.21E+13	https://www.swiggy.com/restaurants/janta-sweet-house-1	Janta Sweet House, Bazar No.9, Circular Road, ...	Menu/531342.json
2	156203	theka coffee desi	Abdohar	3.8	100+ ratings	100.0	2.21E+13	https://www.swiggy.com/restaurants/theka-coffee-desi-1	theka coffee desi, sahitya sadan road city	Menu/156203.json
3	187922	Singh Hut	Abdohar	3.7	20+ ratings	250.0	2.21E+13	https://www.swiggy.com/restaurants/singh-hut-1	Singh Hut, CIRCULAR ROAD NEAR NEHRU PARK, ABDOHAR	Menu/187922.json
4	545330	GRILL MASTERS	Abdohar	4.0	Too Few Ratings	250.0	1.21E+13	https://www.swiggy.com/restaurants/grill-masters-1	GRILL MASTERS, ADA HEIGHTS, Abdohar - Hanumanaga...	Menu/545330.json

Display information of every column(non null values count , data type)

```
In [4]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 148541 entries, 0 to 148540
Data columns (total 10 columns):
 0  Column  Non-Null Count  Dtype
--  --
 1  name      148541 non-null    object
 2  city      148541 non-null    object
 3  rating    148455 non-null    object
 4  rating_count  148455 non-null  object
 5  cost      148419 non-null    float64
 6  lic_no    148512 non-null    object
 7  link      148541 non-null    object
 8  address   148455 non-null    object
 9  menu      148541 non-null    object
dtypes: float64(1), int64(1), object(8)
memory usage: 11.3+ MB
```

Display null values

```
In [5]: df.isnull().sum()
```

```
Out[5]: id                0
name                86
city                9
rating             86
rating_count       86
cost              133
lic_no             29
link               9
address            86
menu              10
dtype: int64
```

Display null values in form of Heatmap

```
In [6]: df.isnull().sum()
```

```
Out[6]: <Axes: >
```

In the dataset, lots of null values are present in every column. Need to handle all null values and clean the dataset.

Approach used to handle null values-

1. Sample values from dataset.
2. Replacing null values with a new category (categorical).
3. Mean replacement.
4. Drop column.

Handle null values

```
In [7]: df['rating'] = pd.to_numeric(df['rating'], replace('-',0), errors='coerce')

In [8]: df['cost'] = pd.to_numeric(df['cost'], replace('nan',0.0), errors='coerce')
df['cost'] = df['cost'].fillna('no')
df['cost'] = df['cost'].astype(float)

In [9]: df['rating'] = df['rating'].fillna(0)
df['rating_count'] = df['rating_count'].replace('-',0), errors='coerce')

In [10]: df['rating_count'] = df['rating_count'].fillna('No Rating')

In [11]: df['address'] = df['address'].fillna('No Address')

In [12]: df['name'] = df['name'].fillna('No')

In [13]: df['lic_no'] = df['lic_no'].str.replace('license','No License')
df['lic_no'] = df['lic_no'].str.replace('license','No License')
df['lic_no'] = df['lic_no'].str.replace('NA','No License')
df['lic_no'] = df['lic_no'].replace(0, 'No License')
df['lic_no'] = df['lic_no'].fillna('No License')

In [14]: df['lic_no'].unique()

Out[14]: array(['2.21E+13', '1.31E+13', 'No License', '1.02E+13', '1.98E+13',
       '2.38E+13', '1.31E+13', '2.11E+13', '1.31E+13', '1.12E+13',
       '2.01E+13', '1.01E+13', '2.25E+13', '1.25E+13', '2.02E+13', '2.02E+13',
       '1.27E+13', '2.27E+13', '2.15E+13', '1.08E+13', '1.07E+13',
       '2.07E+13', '8.48E+13', '4.44E+13', '2.06E+13', '2.90E+13',
       '2.08E+13', '1.14E+13', '1.08E+13', '0.00E+13', '1.22E+13',
       '1.15E+13', '2.08E+13', '2.18E+13', '1.18E+13', '2.22E+13',
       '2.13E+13', '1.13E+13', '2.28E+13', '1.08E+13', '2.08E+13',
       '1.05E+13', '1.24E+13', '2.24E+13', 'ACAG515H464261', '1.04E+13',
       '2.04E+13', '0', '1.09E+13', '2.09E+13', '2.12E+13', '2.14E+13',
       '2.08E+13', '1.28E+13', '2.38E+13', '9.03E+13', '1.38E+13',
       '9.07E+13', '1.72E+13', '1.19E+13', '9.11E+13', '6.48E+13',
       '1.22E+13', '5.21E+13', '1.26E+13', '1.00E+13', '1.55E+13',
       '2.03E+13', '1.03E+13', '3.07E+13', '1.23E+13', '1.06E+13',
       '2.32E+13', '2.38E+13', '1.42E+13', '2.20E+13', '2.24E+12',
       '9.54E+13', '1.08E+13', '1.38E+13', '2.32E+13', '1333900000579',
       '9.61E+13', '2.08E+13', '9.08E+13', '2.07E+13', '2.23E+13',
       '2.19E+13', '1.33E+13', '1.08E+13', '9.02E+13', '2.33E+13',
       '1.16E+13', '2.16E+13', '1.02E+13', '2.10E+13', '1.10E+13',
       '2.09E+13', '2.29E+13', '1.01E+12', '4.02E+13', '2.40E+13',
       '2.38E+13', '9.06E+13', '2.52E+13', '2.77E+13', '1.36E+12',
       '2.35E+13', '7.27E+13', '8C18080121020', '2.29E+13', '1.34E+13',
       '1.22E+13', '4.21E+13', '2.04E+13', '1.27E+13', '2.17E+13',
       '4.21E+13', '4.09E+13', '9.12E+13', '9.83E+13', '3.83E+13'],
      dtype=object)
```

Check null values after handle null values . All null values are repalced and dataset is clean.

```
In [5]: df.isnull().sum()
```

```
Out[5]: id                0
name                0
city                0
rating             0
rating_count       0
cost              0
lic_no             0
link              0
address            0
menu              0
dtype: int64
```

```
In [16]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 148541 entries, 0 to 148540
Data columns (total 10 columns):
 0  Column  Non-Null Count  Dtype
--  --
 1  name      148541 non-null    object
 2  city      148541 non-null    object
 3  rating    148541 non-null    float64
 4  rating_count  148541 non-null    object
 5  cost      148541 non-null    float64
 6  lic_no    148541 non-null    object
 7  link      148541 non-null    object
 8  address   148541 non-null    object
 9  menu      148541 non-null    object
dtypes: float64(2), int64(1), object(7)
memory usage: 11.3+ MB
```

Display Heatmap

```
In [17]: sns.heatmap(df.isnull())
```

```
Out[17]: <Axes: >
```

Display descriptive statistics of dataset

```
In [18]: df.head()
```

	id	name	city	rating	rating_count	cost	lic_no	link	address	menu
0	567335	AB FOODS POINT	Abdohar	4.0	Too Few Ratings	200.0	2.21E+13	https://www.swiggy.com/restaurants/ab-foods-point-1	AB FOODS POINT, NEAR RISH NARANG DENTAL CLINIC, ...	Menu/567335.json
1	531342	Janta Sweet House	Abdohar	4.4	50+ ratings	200.0	1.21E+13	https://www.swiggy.com/restaurants/janta-sweet-house-1	Janta Sweet House, Bazar No.9, Circular Road, ...	Menu/531342.json
2	156203	theka coffee desi	Abdohar	3.8	100+ ratings	100.0	2.21E+13	https://www.swiggy.com/restaurants/theka-coffee-desi-1	theka coffee desi, sahitya sadan road city	Menu/156203.json
3	187922	Singh Hut	Abdohar	3.7	20+ ratings	250.0	2.21E+13	https://www.swiggy.com/restaurants/singh-hut-1	Singh Hut, CIRCULAR ROAD NEAR NEHRU PARK, ABDOHAR	Menu/187922.json
4	545330	GRILL MASTERS	Abdohar	4.0	Too Few Ratings	250.0	1.21E+13	https://www.swiggy.com/restaurants/grill-maste...	GRILL MASTERS, ADA HEIGHTS, Abdohar - Hanumanaga...	Menu/545330.json

Outlier

```
In [19]: sns.boxplot(data=df, x='cost', color='blue')
```

```
Out[19]: <Axes: xlabel='cost'>
```



Outlier

```
In [20]: sns.distplot(df['cost'])
```

```
C:\Users\hp\AppData\Local\Temp\ipykernel_4388\1352492925.py:1: UserWarning:
'distplot' is a deprecated function and will be removed in seaborn v0.14.0.
Please adapt your code to use either 'displot' (a figure-level function with
similar flexibility) or 'histplot' (an axes-level function for histograms).
For a guide to updating your code to use the new functions, please see
https://gist.github.com/mwaskom/d64147ed297445a06372580b65751

sns.distplot(df['cost'])
C:\Users\hp\Anaconda3\Lib\site-packages\seaborn\_oldcore.py:1119: FutureWarning: use_inf_as_na option is deprecated and will be removed in a future version. Convert inf values to NaN before operating instead.
  with pd.option_context('mode.use_inf_as_na', True):
  with pd.option_context('mode.use_inf_as_na', True):
<Axes: xlabel='cost', ylabel='Density'>
```



Approach Use-

Remove outlier with IQR (Interquartile Range)

```
In [21]: std=df['cost'].std()
```

```
Out[21]: 796.4258426177296
```

```
In [22]: iqr=df['cost'].quantile(0.75)-df['cost'].quantile(0.25)
```

```
Out[22]: 180.8
```

```
In [23]: upperlimit=iqr*1.5*std
```

```
lowerlimit=iqr*1.5*std
```

```
In [24]: upperlimit
```

```
Out[24]: 1284.6375639265943
```

```
In [25]: lowerlimit
```

```
Out[25]: -1894.6375639265943
```

```
In [26]: df['cost'] = df[(df['cost'] >= lowerlimit) & (df['cost'] <= upperlimit)]['cost']
```

```
In [27]: std=df['cost'].std()
```

```
Out[27]: 145.56961618948744
```

```
In [28]: iqr=df['cost'].quantile(0.75)-df['cost'].quantile(0.25)
```

```
Out[28]: 180.8
```

```
In [29]: upperlimit=iqr*1.5*std
```

```
lowerlimit=iqr*1.5*std
```

```
In [30]: upperlimit
```

```
Out[30]: 318.35442418423115
```

```
In [31]: lowerlimit
```

```
Out[31]: -118.35442418423115
```

```
In [32]: df.loc[df['cost']>318,'cost']=318
```

```
In [33]: lowerlimit = -118
```

```
df['cost'] = df[(df['cost'] >= lowerlimit)]['cost']
```

```
In [34]: std=df['cost'].std()
```

```
Out[34]: 67.9888890287151
```

```
In [35]: iqr=df['cost'].quantile(0.75)-df['cost'].quantile(0.25)
```

```
Out[35]: 180.8
```

```
In [36]: upperlimit=iqr*1.5*std
```

```
lowerlimit=iqr*1.5*std
```

```
In [37]: upperlimit
```

```
Out[37]: -1.983333543072627
```

```
In [38]: df['cost'] = df[(df['cost'] >= lowerlimit)]['cost']
```

```
In [39]: std=df['cost'].std()
```

```
Out[39]: 67.9888890287151
```

```
In [40]: df = df[(df['cost'] >= max(0, lowerlimit))]
```

```
Out[40]: lowerlimit
```

```
Out[40]: -1.983333543072627
```

```
In [41]: lowerlimit = 0
```

```
df['cost'] = df[(df['cost'] >= lowerlimit)]['cost']
std_dev = df['cost'].std()
print(std_dev)
```

```
df = df[(df['cost'] >= max(0, lowerlimit))]
```

```
Out[41]: 67.9888890287151
```

```
In [42]: lowerlimit
```

```
Out[42]: 0
```

```
In [43]: std_dev = df['cost'].std()
```

```
Out[43]: 67.9888890287151
```

```
In [44]: df = df[(df['cost'] >= max(0, lowerlimit))]
```

```
Out[44]: lowerlimit
```

```
Out[44]: -1.983333543072627
```

```
In [45]: lowerlimit = 0
```

```
df['cost'] = df[(df['cost'] >= lowerlimit)]['cost']
std_dev = df['cost'].std()
print(std_dev)
```

```
df = df[(df['cost'] >= max(0, lowerlimit))]
```

```
Out[45]: 67.9888890287151
```

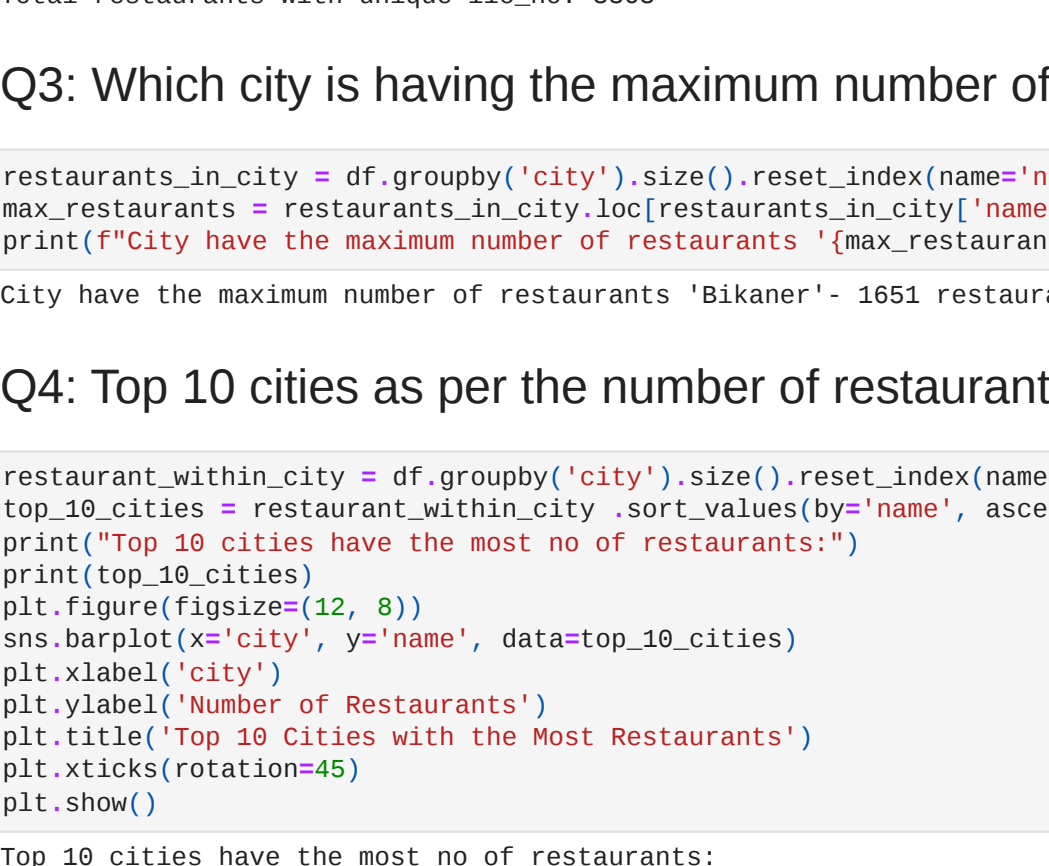
```
In [46]: df = df[(df['cost'] >= max(0, lowerlimit))]
```

```
Out[46]: lowerlimit
```

```
Out[46]: 0
```

```
In [51]: sns.boxplot(data=df, x='cost', color='red')
```

```
Out[51]: <Axes: xlabel='cost'>
```



```
Q1 = df['cost'].quantile(0.25)
```

```
Q3 = df['cost'].quantile(0.75)
```

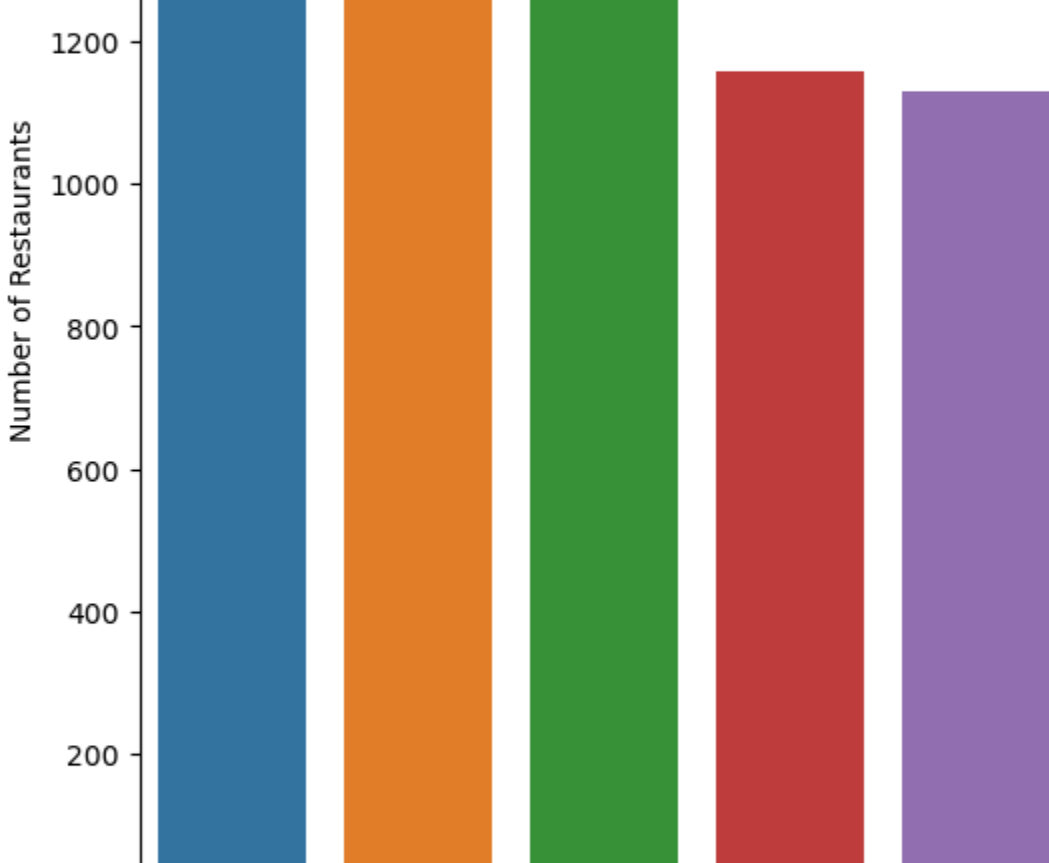
```
IQR = Q3 - Q1
```

```
filter = (df['cost'] >= Q1 - 1.5 * IQR) & (df['cost'] <= Q3 + 1.5 * IQR)
```

```
df = df.loc[filter]
```

```
In [53]: sns.boxplot(data=df, x='cost', color='green')
```

```
Out[53]: <Axes: xlabel='cost'>
```

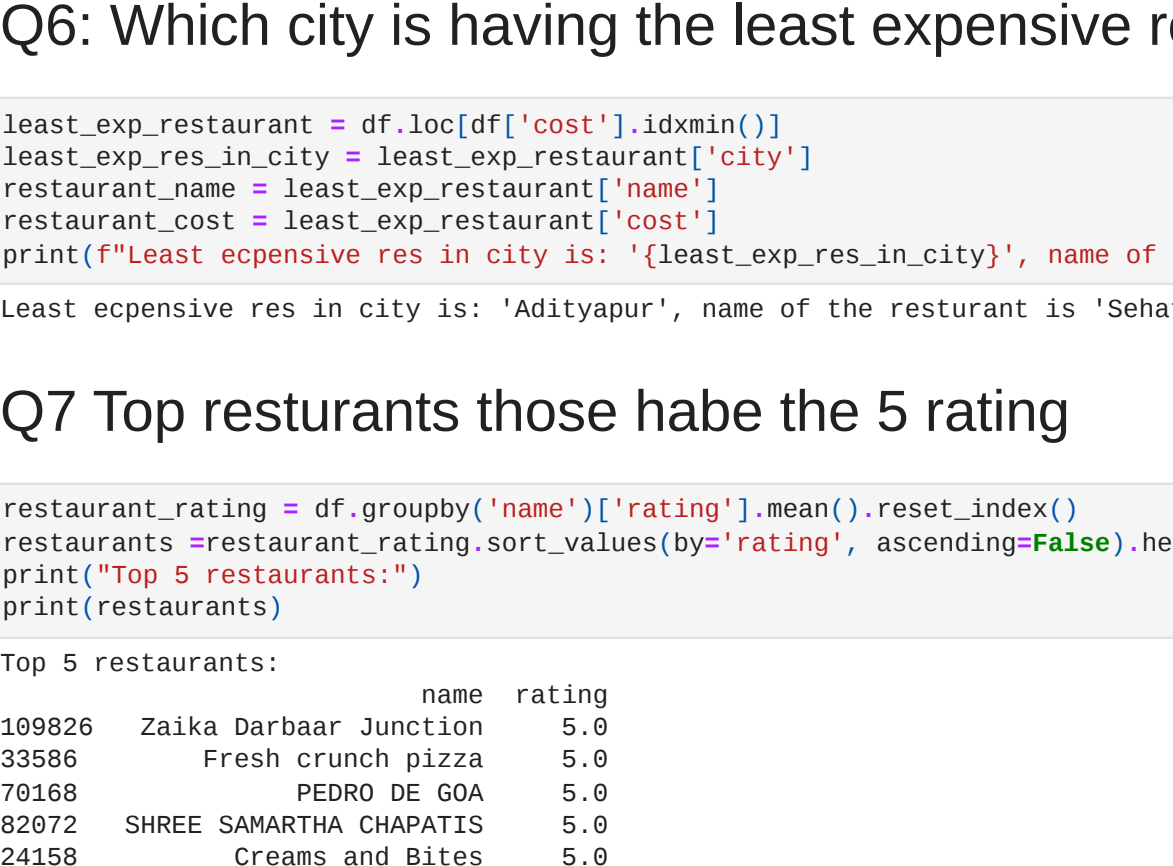


Distribution plot after remove outliers

```
In [54]: sns.distplot(df['cost'])
```

```
C:\Users\hp\AppData\Local\Temp\ipykernel_4388\1352492925.py:1: UserWarning:
'distplot' is a deprecated function and will be removed in seaborn v0.14.0.
Please adapt your code to use either 'displot' (a figure-level function with
similar flexibility) or 'histplot' (an axes-level function for histograms).
For a guide to updating your code to use the new functions, please see
https://gist.github.com/mwaskom/d64147ed297445a06372580b65751

sns.distplot(df['cost'])
C:\Users\hp\Anaconda3\Lib\site-packages\seaborn\_oldcore.py:1119: FutureWarning: use_inf_as_na option is deprecated and will be removed in a future version. Convert inf values to NaN before operating instead.
  with pd.option_context('mode.use_inf_as_na', True):
  with pd.option_context('mode.use_inf_as_na', True):
<Axes: xlabel='cost', ylabel='Density'>
```



Q1: How many restaurants are listed on Swiggy?

```
In [55]: unique_restaurants = df['name'].nunique()
```

```
print(f'Total number of restaurants listed on Swiggy: {unique_restaurants}')

Total number of restaurants listed on Swiggy: 12129
```

Q2: How many cities are having restaurants that are listed on Swiggy?

```
In [56]: df.lic_no = df[df['lic_no'] != notna]]
```

```
list_restaurants = df.lic_no.groupby('lic_no')['city'].nunique().reset_index()
```

```
print(f'Top 10 cities have the most no of restaurants')
restaurants_lic_no = list_restaurants.sort_values('city', ascending=False).head(10)
```

```
print(f'Total restaurants with unique lic_no: {restaurants_lic_no}')

Total restaurants with unique lic_no: 3393
```

Q3: Which city is having the maximum number of restaurants?

```
In [57]: restaurants_in_city = df.groupby('city').size().reset_index(name='count')
```

```
max_restaurants = restaurants_in_city.loc[restaurants_in_city['name'].idxmax()]
```

```
print(f'The city having the maximum number of restaurants is {max_restaurants["city"]} with rating of {max_restaurants["rating"]}')

The city having the maximum number of restaurants is 'Bikarner' - 1651 restaurants.
```

Q4: Top 10 cities as per the number of restaurants listed?

```
In [58]: restaurant_within_city = df.groupby('city').size().reset_index(name='count')
```

```
top_10_cities = restaurant_within_city.sort_values('count', ascending=False).head(10)
```

```
print(f'Top 10 cities have the most no of restaurants')
plt.figure(figsize=(12,8))
```

```
sns.barplot(x='city', y='count', data=top_10_cities)
```

```
plt.xlabel('City')
```

```
plt.ylabel('Number of Restaurants')
```

```
plt.title('Top 10 Cities with the Most Restaurants')
```

```
plt.xticks(rotation=45)
```

```
plt.show()
```

Top 10 cities have the most no of restaurants:

city	count
Bikarner	1651
Noida	1410
Indraprastha	1278
BTH, Bangalore	1168
Rohtak	1130
Kothrud, Pune	1084
Indiranagar, Bangalore	1061
Electronic City, Bangalore	1037
Greater Kailash 2, Delhi	1027
Vashi, Mumbai	1019

With the help of this bar graph, we can easily see the top 10 cities where Swiggy has the most restaurants present.

Q5: Most Popular hotel by rating throughout the dataset?

```
In [59]: hotel_rating = df.groupby('name')['rating'].mean().reset_index()
```