

# Restaurant Recommender

In this project we are going to assist a restaurant consolidator to build a restaurant recommender which will help a customer find the best restaurants in an area and find out which cuisines are served best so on and so forth. Factors such as delivery options, ratings and price options will be explored and the insights will be presented. Also we will find out what makes a restaurant more appealing to the customer and which factors make the ratings go up or down.

## Step 1: Setting up the work environment

We are going to download the necessary packages for our work. We are going to view the dataset and check the datatypes.

```
In [1]: #installing packages
import numpy as np
import pandas as pd
import matplotlib as plt
import seaborn as sb
from scipy import stats
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error
from sklearn.pipeline import Pipeline
from sklearn.preprocessing import StandardScaler
from sklearn.preprocessing import PolynomialFeatures
from sklearn.metrics import r2_score
from sklearn.model_selection import train_test_split
from sklearn.model_selection import cross_val_score
from sklearn.model_selection import cross_val_predict
from sklearn.linear_model import Ridge
from sklearn.model_selection import GridSearchCV
import os
from collections import Counter
from numpy import mean
from numpy import std
from pandas import read_csv

from sklearn.dummy import DummyClassifier,
from sklearn.preprocessing import StandardScaler
from sklearn.neighbors import KNeighborsClassifier
from sklearn.pipeline import Pipeline
from sklearn.model_selection import cross_val_score
from sklearn.metrics import RepeatedStratifiedKFold
from sklearn.metrics import precision_recall_curve
from sklearn.metrics import auc
from sklearn.metrics import make_scorer
from sklearn.tree import DecisionTreeClassifier
from sklearn.neighbors import KNeighborsClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.ensemble import ExtraTreesClassifier
from sklearn.ensemble import BaggingClassifier
from sklearn.naive_bayes import GaussianNB
from pandas import Series
from numpy.random import randn
from statsmodels.stats.weightstats import _test_ind
import scipy.stats as stats
from scipy.stats import ttest_ind

%matplotlib inline
from matplotlib import pyplot as plt
from scipy.stats import shapiro
from scipy import stats
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error
from sklearn.metrics import classification_report
from sklearn.pipeline import Pipeline
from sklearn.preprocessing import StandardScaler
from sklearn.preprocessing import PolynomialFeatures
from sklearn.metrics import r2_score
from sklearn.model_selection import train_test_split
from sklearn.model_selection import cross_val_score
from sklearn.model_selection import cross_val_predict
from sklearn.linear_model import Ridge
from sklearn.model_selection import GridSearchCV
from datetime import datetime
from datetime import date
import statsmodels.api as sm
from statsmodels.formula.api import ols
from sklearn.preprocessing import StandardScaler
from sklearn.decomposition import PCA
from sklearn.linear_model import LogisticRegression as lgr

from sklearn import preprocessing as preproc
from sklearn import metrics
from sklearn.datasets import make_classification
from sklearn.model_selection import train_test_split
from sklearn.pipeline import make_pipeline
from sklearn.preprocessing import StandardScaler
from sklearn.neighbors import KNeighborsRegressor
from sklearn.neural_network import MLPRegressor

#downloading the file
cmd = pd.read_excel("C:\Users\aujoydutta\Desktop\Data analysis\Datasets for ML\Regression\data.xlsx")
rr.head()
```

	Restaurant ID	Restaurant Name	Country Code	City	Address	Locality	Locality Verbose	Longitude	Latitude	Cuisines	cuisine count	Average Cost for two	Currency	Tab booking
0	53	Amber	1	New Delhi	N-19, Connaught Place, New Delhi	Connaught Place	Connaught Place, New Delhi	77.220891	28.630197	North Indian, Chinese, Mughlai	3.0	1800	Indian Rupees(INR)	Yes
1	55	Berco's	1	New Delhi	G-2/43, Middle Circle Connaught Place, New Delhi	Connaught Place	Connaught Place, New Delhi	77.217298	28.632452	Chinese, Thai	2.0	1100	Indian Rupees(INR)	Yes
2	60	Colonel's Kababz	1	New Delhi	29, Defence Colony Market, Defence Colony, New...	Defence Colony	Defence Colony, New Delhi	77.230591	28.574036	North Indian, Mughlai	2.0	900	Indian Rupees(INR)	Yes
3	64	Diva -The Italian Restaurant	1	New Delhi	M-8A, M Block Market, Greater Kailash (GK 2), ...	Greater Kailash (GK 2)	Greater Kailash (GK 2), New Delhi	77.243186	28.534202	Italian	1.0	2500	Indian Rupees(INR)	Yes
4	65	Drums of Heaven	1	New Delhi	S-14, Green Park Extension, Green Park, New Delhi	Green Park	Green Park, New Delhi	77.205934	28.558018	Chinese, Seafood, Thai	3.0	1800	Indian Rupees(INR)	Yes

```
In [3]: #getting the secondary data
cmd = pd.read_excel("C:\Users\aujoydutta\Desktop\Data analysis\Datasets for DV\Country-Code.xlsx")
rr.head()
```

	Country Code	Country
0	1	India
1	14	Australia
2	30	Brazil
3	37	Canada
4	94	Indonesia

```
In [4]: #merging secondary into primary
cmd.merge(cmd, on="Country Code", how="left")
rr.head()
```

	Restaurant ID	Restaurant Name	Country Code	City	Address	Locality	Locality Verbose	Longitude	Latitude	Cuisines	...	Average Cost for two	Currency	Tab booking
0	53	Amber	1	New Delhi	N-19, Connaught Place, New Delhi	Connaught Place	Connaught Place, New Delhi	77.220891	28.630197	North Indian, Chinese, Mughlai	...	1800	Indian Rupees(INR)	Yes
1	55	Berco's	1	New Delhi	G-2/43, Middle Circle Connaught Place, New Delhi	Connaught Place	Connaught Place, New Delhi	77.217298	28.632452	Chinese, Thai	...	1100	Indian Rupees(INR)	Yes
2	60	Colonel's Kababz	1	New Delhi	29, Defence Colony Market, Defence Colony, New...	Defence Colony	Defence Colony, New Delhi	77.230591	28.574036	North Indian, Mughlai	...	900	Indian Rupees(INR)	Yes
3	64	Diva -The Italian Restaurant	1	New Delhi	M-8A, M Block Market, Greater Kailash (GK 2), ...	Greater Kailash (GK 2)	Greater Kailash (GK 2), New Delhi	77.243186	28.534202	Italian	...	2500	Indian Rupees(INR)	Yes
4	65	Drums of Heaven	1	New Delhi	S-14, Green Park Extension, Green Park, New Delhi	Green Park	Green Park, New Delhi	77.205934	28.558018	Chinese, Seafood, Thai	...	1800	Indian Rupees(INR)	Yes

5 rows x 21 columns

```
In [5]: #examining the dataset
rr.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 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     9550 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   cuisine count      9550 non-null   float64
11   Average Cost for two 9551 non-null   int64
12   Currency           9551 non-null   int64
13   Has Table booking   9551 non-null   object
14   Has online delivery 9551 non-null   object
15   Price range        9551 non-null   object
16   Aggregate rating    9551 non-null   float64
17   Rating color       9551 non-null   object
18   Rating text        9551 non-null   object
19   Votes              9551 non-null   int64
20   Country            9551 non-null   object
dtypes: float64(4), int64(5), object(12)
memory usage: 1.6+ MB
```

```
In [6]: #examining the dataset
rr.shape

(9551, 21)
```

## Step 2: Data cleaning

In this step, we are going to clean our dataset. We are going to look for null values and replace them with mean and mode. We are going to modify some variables if it is necessary and change datatypes for better analysis. We will also remove outliers from the dataset. Outliers hamper the machine learning algorithms and hence they have to be removed.

```
In [7]: # column name cleaning
rr.columns = rr.columns.str.replace(' ', '')
rr.columns

Index(['RestaurantID', 'RestaurantName', 'CountryCode', 'City', 'Address',
       'Locality', 'LocalityVerbose', 'Longitude', 'Latitude', 'Cuisines',
       'cuisinecount', 'AverageCostfortwo', 'Currency', 'HasTablebooking',
       'HasOnlineDelivery', 'Pricerange', 'AggregateRating', 'Ratingcolor',
       'Ratingtext', 'Votes', 'Country'],
      dtype='object')
```

```
In [8]: #checking null values
rr.isna().sum(axis=0)

RestaurantID      0
RestaurantName    1
CountryCode       0
City              0
Address           0
Locality          0
LocalityVerbose  0
Longitude         0
Latitude          0
Cuisines          9
cuisinecount     1
AverageCostfortwo 0
Currency          0
HasTablebooking  0
HasOnlineDelivery 0
Pricerange        0
AggregateRating   0
Ratingcolor       0
Ratingtext        0
Votes             0
Country           0
dtype: int64
```

```
In [9]: #drop null
rr=rr.dropna()
```

	RestaurantID	RestaurantName	CountryCode	City	Address	Locality	LocalityVerbose	Longitude	Latitude	Cuisines	...	Average Cost for two	Currency	Tab booking
0	53	Amber	1	New Delhi	N-19, Connaught Place, New Delhi	Connaught Place	Connaught Place, New Delhi	77.220891	28.630197	North Indian, Chinese, Mughlai	...	1800	Indian Rupees(INR)	Yes
1	55	Berco's	1	New Delhi	G-2/43, Middle Circle Connaught Place, New Delhi	Connaught Place	Connaught Place, New Delhi	77.217298	28.632452	Chinese, Thai	...	1100	Indian Rupees(INR)	Yes
2	60	Colonel's Kababz	1	New Delhi	29, Defence Colony Market, Defence Colony, New...	Defence Colony	Defence Colony, New Delhi	77.230591	28.574036	North Indian, Mughlai	...	900	Indian Rupees(INR)	Yes
3	64	Diva -The Italian Restaurant	1	New Delhi	M-8A, M Block Market, Greater Kailash (GK 2), ...	Greater Kailash (GK 2)	Greater Kailash (GK 2), New Delhi	77.243186	28.534202	Italian	...	2500	Indian Rupees(INR)	Yes
4	65	Drums of Heaven	1	New Delhi	S-14, Green Park Extension, Green Park, New Delhi	Green Park	Green Park, New Delhi	77.205934	28.558018	Chinese, Seafood, Thai	...	1800	Indian Rupees(INR)	Yes

9540 rows x 21 columns

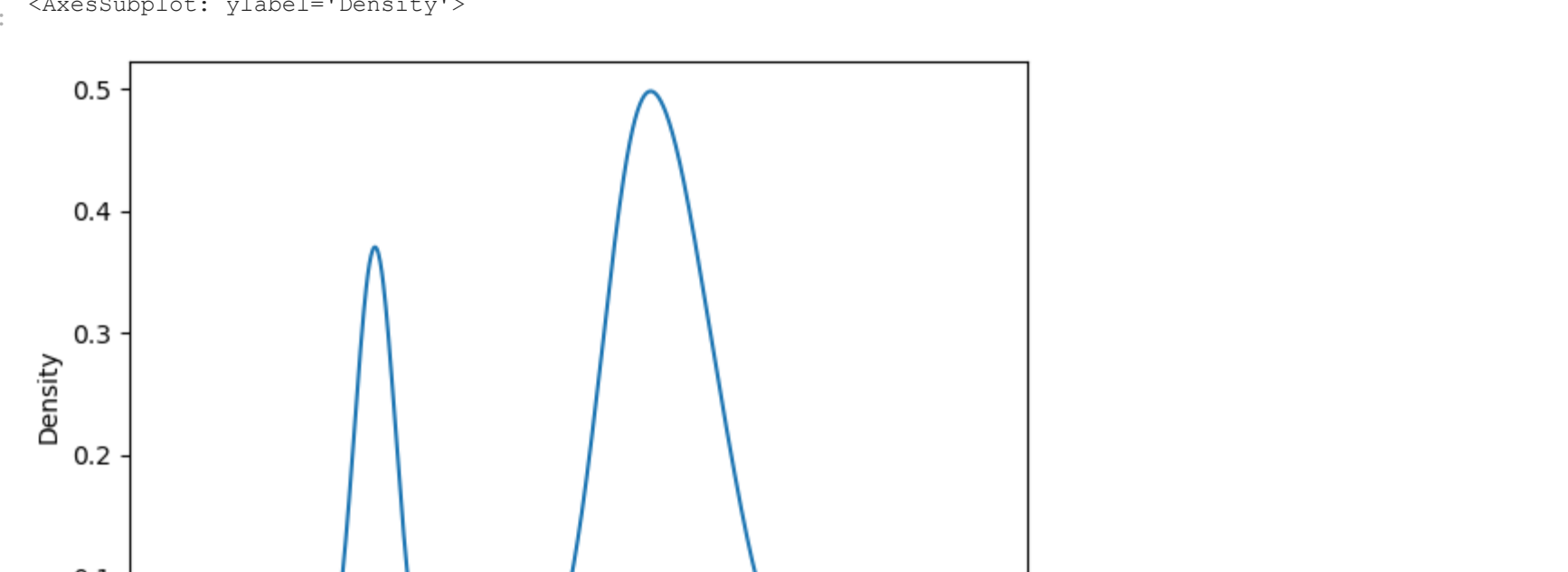
```
In [10]: #dropping useless columns
rr=rr.drop(['Address', 'Locality', 'LocalityVerbose', 'Longitude', 'Latitude', 'Ratingcolor', 'RestaurantID'], axis=1)
rr.head()
```

	RestaurantName	CountryCode	City	Cuisines	cuisinecount	AverageCostfortwo	Currency	HasTableBooking	HasOnlineDelivery	Pricerange
0	Amber	1	New Delhi	North Indian, Chinese, Mughlai	3.0	1800	Indian Rupees(INR)	Yes	Yes	
1	Berco's	1	New Delhi	Chinese, Thai	2.0	1100	Indian Rupees(INR)	Yes	Yes	
2	Colonel's Kababz	1	New Delhi	North Indian, Mughlai	2.0	900	Indian Rupees(INR)	Yes	No	
3	Diva -The Italian Restaurant	1	New Delhi	Italian	1.0	2500	Indian Rupees(INR)	Yes	Yes	
4	Drums of Heaven	1	New Delhi	Chinese, Seafood, Thai	3.0	1800	Indian Rupees(INR)	Yes	Yes	

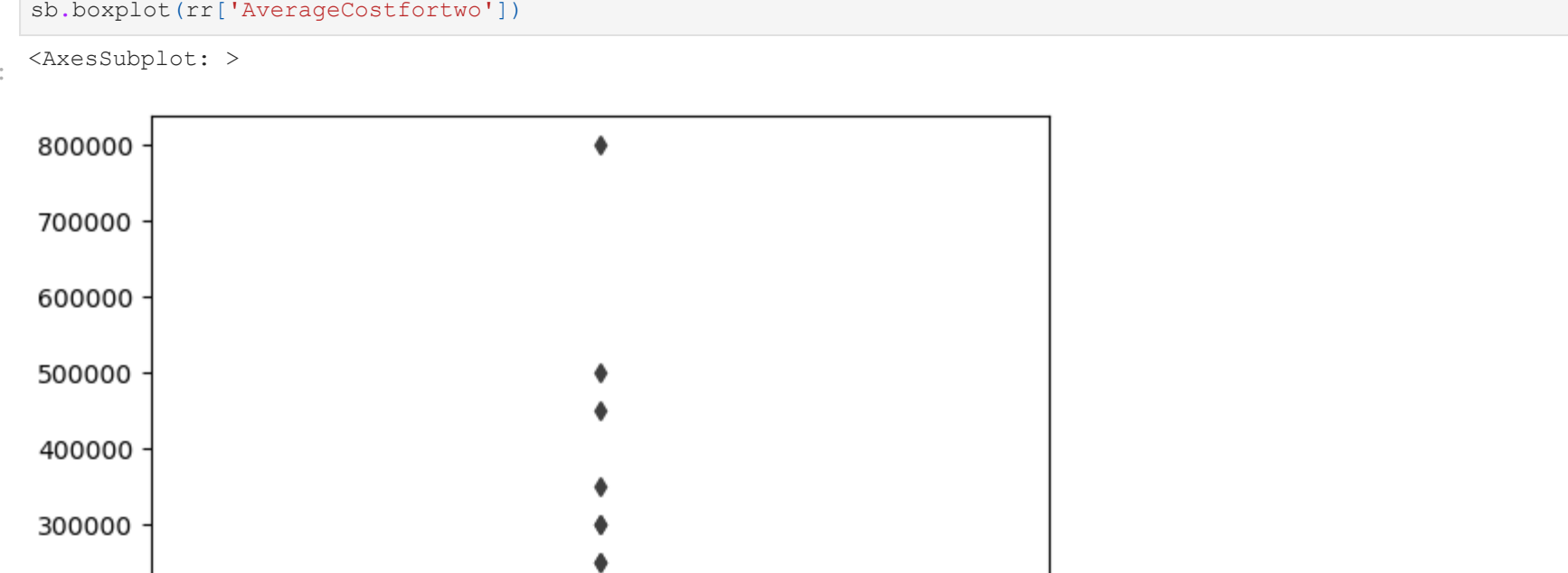
## Univariate analysis

We are using boxplot, histogram and plot to find the outliers, extreme values and distribution of the numerical variables.

```
In [11]: #checking distribution for Average Cost for two
rr['AverageCostfortwo'].plot.density()
```



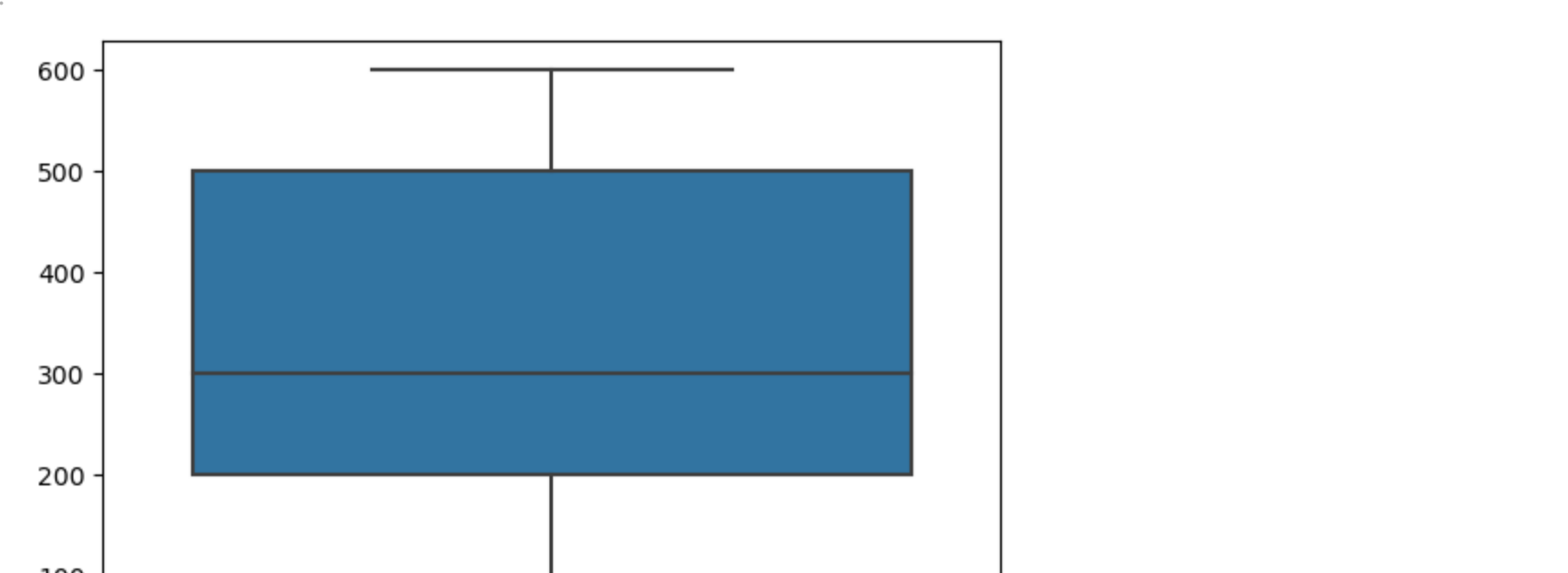
```
In [12]: #checking distribution for Votes
rr['Votes'].plot.density()
```



```
In [13]: #checking distribution for Aggregate rating
rr['AggregateRating'].plot.density()
```



```
In [14]: #checking outliers using boxplot Average Cost for two
sb.boxplot(rr['AverageCostfortwo'])
```



```
In [15]: #outlier treatment for variable Average Cost for two
rr['AverageCostfortwo'].quantile([0.1, 0.25, 0.5, 0.7, 0.9, 0.95, 0.99])
```

Out[15]:

0.10	100.0
0.25	250.0
0.50	400.0
0.70	600.0
0.90	1200.0
0.95	1700.0
0.99	3300.0
Name:	AverageCostfortwo, dtype: float64

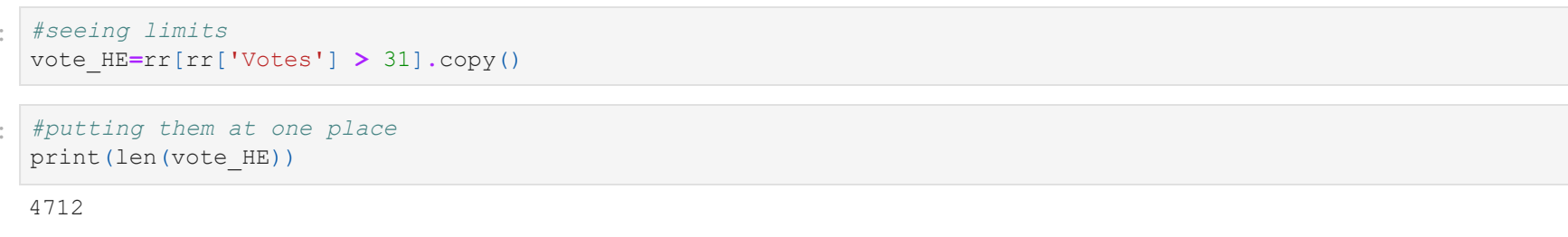
```
In [16]: #setting limits
AverageCosttwo_RE=rr[rr['AverageCostfortwo'] > 600].copy()
```

```
In [17]: #putting them at one place
print(len(AverageCosttwo_RE))

2636
```

```
In [18]: #removing outliers
for x in rr['AverageCostfortwo']:
    if x > 600.0:
        rr['AverageCostfortwo'].replace(x,np.nan,inplace=True)
```

```
In [19]: #checking outliers gone or not
sb.boxplot(rr['AverageCostfortwo'])
```



```
In [20]: #replacing with median values
medcost=rr['AverageCostfortwo'].median()
print(medcost)
rr['AverageCostfortwo'].replace(np.nan, medcost, inplace=True)
```

300.0

```
In [21]: #checking outliers using boxplot Votes
sb.boxplot(rr['Votes'])
```



```
In [22]: #outlier treatment for variable votes
rr['Votes'].quantile([0.1, 0.25, 0.5, 0.7, 0.9, 0.95, 0.99])
```

Out[22]:

0.10	0.00
0.25	5.00
0.50	31.00
0.70	97.00
0.90	376.10
0.95	697.00
0.99	1883.49
Name:	Votes, dtype: float64

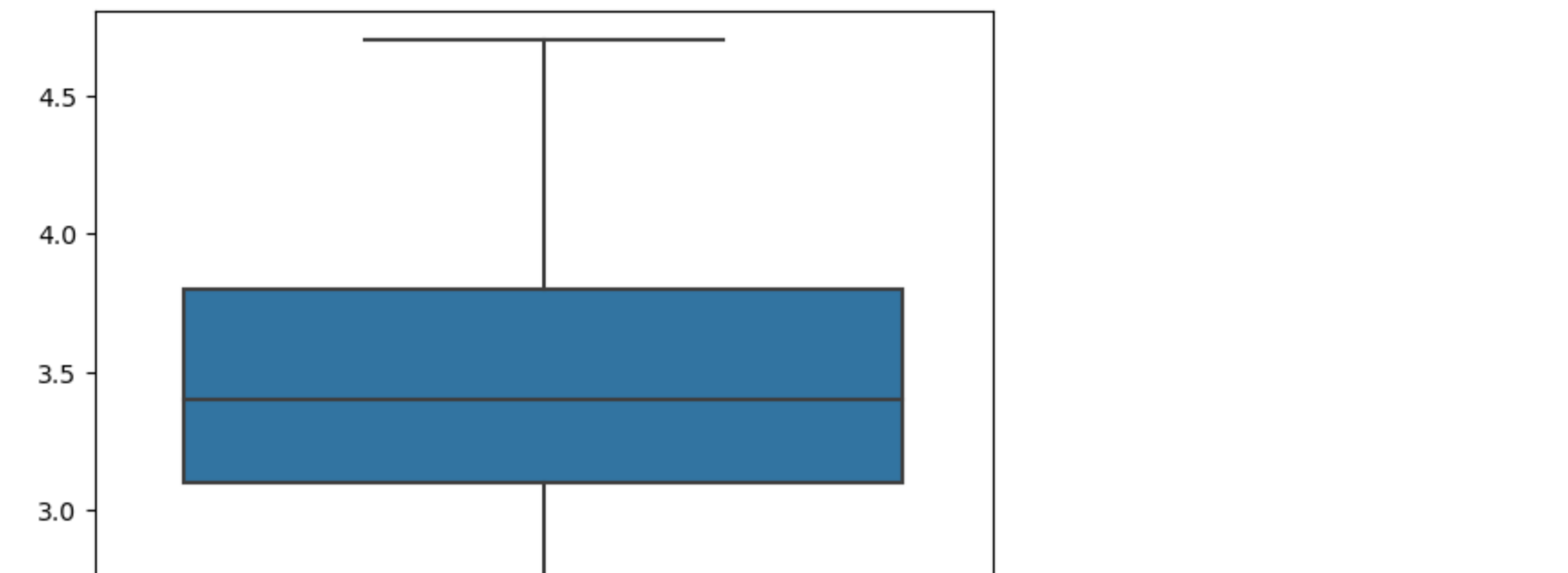
```
In [23]: #setting limits
vote_RE=rr[rr['Votes'] > 31].copy()
```

```
In [24]: #putting them at one place
print(len(vote_RE))

4712
```

```
In [25]: #removing outliers
for x in rr['Votes']:
    if x > 31:
        rr['Votes'].replace(x,np.nan,inplace=True)
```

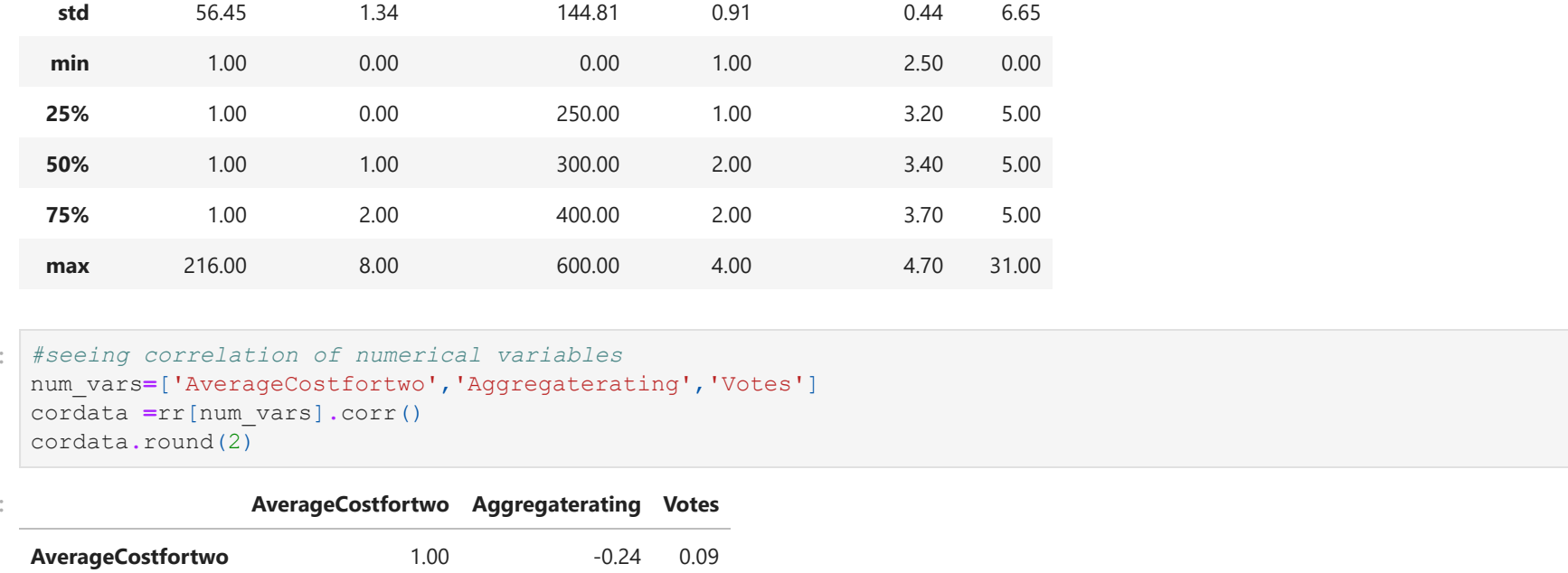
```
In [26]: #checking if outliers gone
sb.boxplot(rr['Votes'])
```



```
In [27]: #replacing with median values
medvote=rr['Votes'].median()
print(medvote)
rr['Votes'].replace(np.nan, medvote, inplace=True)
```

5.0

```
In [28]: #checking outliers using boxplot Aggregate rating
sb.boxplot(rr['AggregateRating'])
```



```
In [29]: #outlier treatment for variable Aggregate rating
rr['AggregateRating'].quantile([0.1, 0.25, 0.5, 0.7, 0.9, 0.95, 0.99])
```

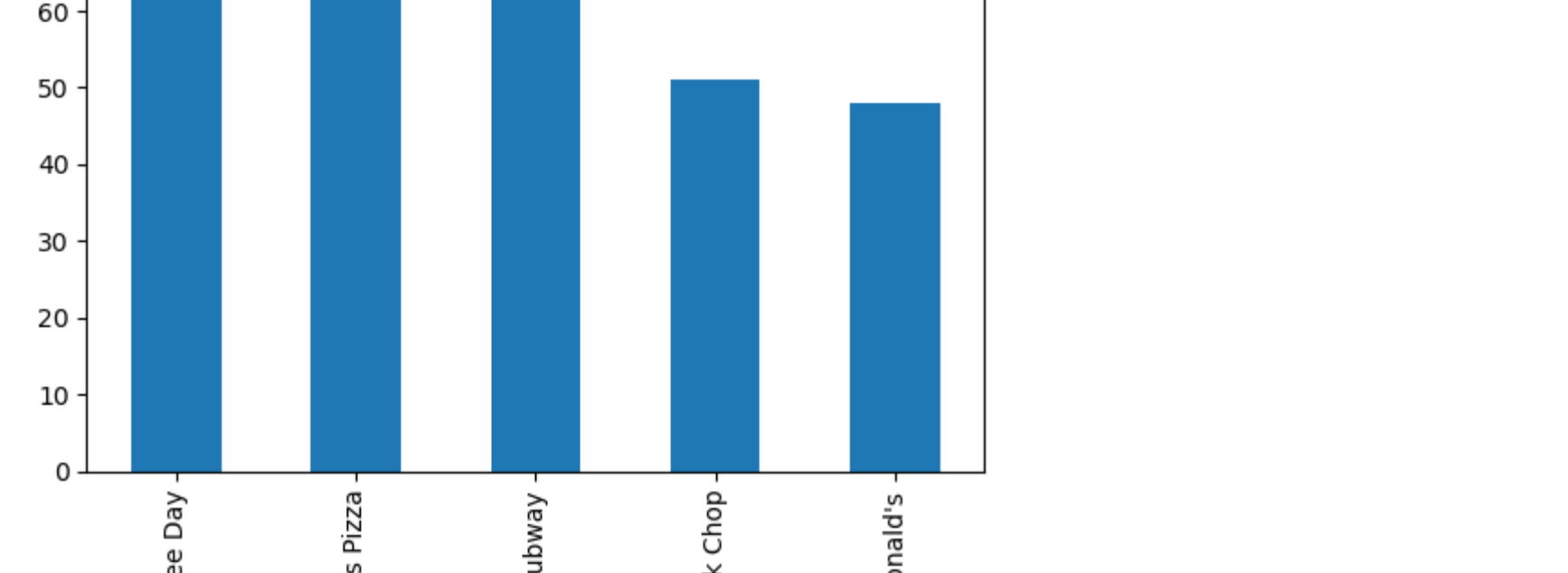
Out[29]:

0.10	0.0
0.25	2.5
0.50	3.2
0.70	3.6
0.90	4.1
0.95	4.3
0.99	4.7
Name:	AggregateRating, dtype: float64

```
In [30]: #setting limits
ar_RE=rr[rr['AggregateRating'] < 2.5].copy()
ar_MED=rr[rr['AggregateRating'] > 4.7].copy()
```

```
In [31]: #removing outliers
for x in rr['AggregateRating']:
    if x < 2.5:
        rr['AggregateRating'].replace(x,np.nan,inplace=True)
for x in rr['AggregateRating']:
    if x > 4.7:
        rr['AggregateRating'].replace(x,np.nan,inplace=True)
```

```
In [32]: #checking if outliers gone
sb.boxplot(rr['AggregateRating'])
```



```
In [33]: #replacing with median values
medar=rr['AggregateRating'].median()
print(medar)
rr['AggregateRating'].replace(np.nan, medar, inplace=True)
```

3.4

## Step 3: Exploratory data analysis

In this step, we are going to explore the dataset. Perform hypothesis tests, bivariate analysis and check for correlation between variables.

Out[34]:

	CountryCode	cuisinecount	AverageCostfortwo	Pricerange	AggregateRating	Votes
count	9540.00	9540.00	9540.00	9540.00	9540.00	9540.00
mean	18.17	1.12	318.56	1.80	3.44	6.73
std	56.45	1.34	144.81	0.91	0.44	6.65
min	1.00	0.00	0.00	1.00	2.50	0.00
25%	1.00	0.00	250.00	1.00	3.20	5.00
50%	1.00	1.00	300.00	2.00	3.40	5.00
75%	1.00	2.00	400.00	2.00	3.70	5.00
max	216.00	8.00	600.00	4.00	4.70	31.00

Out[35]:

	AverageCostfortwo	AggregateRating	Votes
AverageCostfortwo	1.00	-0.24	0.09
AggregateRating	-0.24	1.00	-0.21
Votes	0.09	-0.21	1.00

```
In [36]: #heatmap of correlation
sb.heatmap(cordata, annot=True, cmap="Greens")
```



```
In [37]: #seeing the top franchise
francnt=rr['City'].value_counts().head(5)
francnt.plot(kind="bar")
```



```
In [38]: #top 3 cities with most restaurants
n_by_city=rr.groupby("City")["RestaurantName"].count().nlargest(3)
n_by_city.plot(kind="bar")
```



