

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 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
cm = pd.read_excel("C:\Users\aujoydutta\Desktop\Data analysis\Datasets for ML\Regression\data.xlsx")
cm.head()
```

Restaurant ID	Restaurant Name	Country Code	City	Address	Locality	Locality Verbose	Longitude	Latitude	Cuisines	cuisine count	Average Cost for two	Currency	
0	53	Amber	1	New Delhi	N-19, Connaught Place, New Delhi	Connaught Place, New Delhi	Connaught Place, New Delhi	77.220891	28.630197	North Indian, Chinese, Mughlai	3.0	1800	Indian Rupees(INR)
1	55	Berco's	1	New Delhi	G-2/43, Middle Circle, Connaught Place, New Delhi	Connaught Place, New Delhi	Connaught Place, New Delhi	77.217298	28.632452	Chinese, Thai	2.0	1100	Indian Rupees(INR)
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, Chinese, Mughlai	2.0	900	Indian Rupees(INR)
3	64	Divs - 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)
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)

```
In [2]: #getting the secondary data
cm = pd.read_excel("C:\Users\aujoydutta\Desktop\Data analysis\Datasets for DV\Country-Code.xlsx")
cm.head()
```

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

```
In [3]: #merging secondary into primary
cm = pd.merge(cm, cm, on="Country Code", how="left")
cm.head()
```

Restaurant ID	Restaurant Name	Country Code	City	Address	Locality	Locality Verbose	Longitude	Latitude	Cuisines	...	Average Cost for two	Currency	Has Table booking	
0	53	Amber	1	New Delhi	N-19, Connaught Place, New Delhi	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, 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, Chinese, Mughlai	...	900	Indian Rupees(INR)	Yes
3	64	Divs - 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            9550 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          1
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()
rr
```

RestaurantID	RestaurantName	CountryCode	City	Address	Locality	LocalityVerbose	Longitude	Latitude	Cuisines	...	Average Cost for two	Currency	Has Table booking	
0	53	Amber	1	New Delhi	N-19, Connaught Place, New Delhi	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, 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, Chinese, Mughlai	...	900	Indian Rupees(INR)	Yes
3	64	Divs - 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()
```

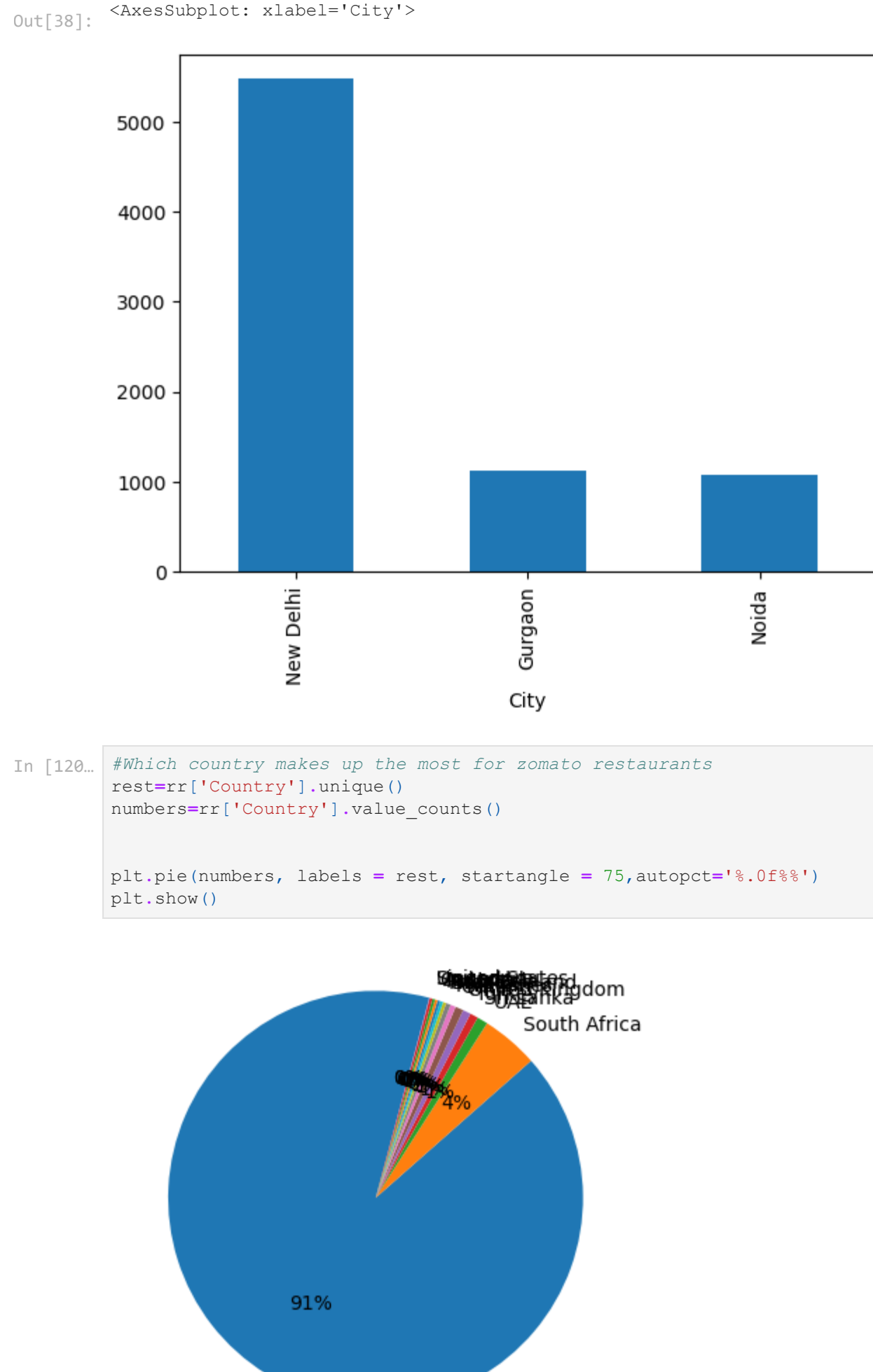
RestaurantName	CountryCode	City	Cuisines	cuisinecount	AverageCostfortwo	Currency	HasTableBooking	HasOnlineDelivery	Pricerange
0	Amber	1	New Delhi, Indian, Chinese, Mughlai	3.0	1800	Indian Rupees(INR)	Yes	Yes	
1	Berco's	1	New Delhi, Chinese, Thai	2.0	1100	Indian Rupees(INR)	Yes	Yes	
2	Colonel's Kababz	1	New Delhi, North Indian, Indian, Mughlai	2.0	900	Indian Rupees(INR)	Yes	No	
3	Divs - 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.

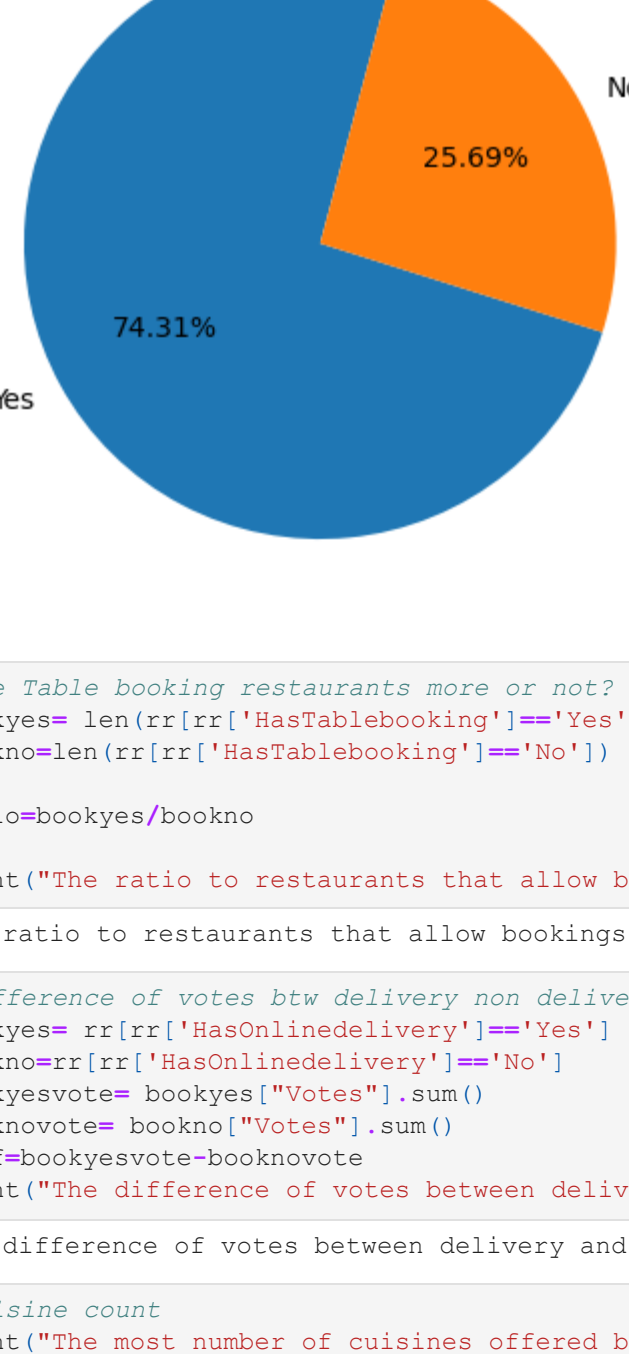
```
In [11]: #checking distribution for Average Cost for two
rr['AverageCostfortwo'].plot.density()
```

AverageCostfortwo	Density
0	0.000000
1	0.000000
2	0.000000
3	0.000000
4	0.000000
5	0.000000
6	0.000000
7	0.000000
8	0.000000
9	0.000000
10	0.000000
11	0.000000
12	0.000000
13	0.000000
14	0.000000
15	0.000000
16	0.000000
17	0.000000
18	0.000000
19	0.000000
20	0.000000
21	0.000000
22	0.000000
23	0.000000
24	0.000000
25	0.000000
26	0.000000
27	0.000000
28	0.000000
29	0.000000
30	0.000000
31	0.000000
32	0.000000
33	0.000000
34	0.000000
35	0.000000
36	0.000000
37	0.000000
38	0.000000
39	0.000000
40	0.000000
41	0.000000
42	0.000000
43	0.000000
44	0.000000
45	0.000000
46	0.000000
47	0.000000
48	0.000000
49	0.000000
50	0.000000
51	0.000000
52	0.000000
53	0.000000
54	0.000000
55	0.000000
56	0.000000
57	0.000000
58	0.000000
59	0.000000
60	0.000000
61	0.000000
62	0.000000
63	0.000000
64	0.000000
65	0.000000
66	0.000000
67	0.000000
68	0.000000
69	0.000000
70	0.000000
71	0.000000
72	0.000000
73	0.000000
74	0.000000
75	0.000000
76	0.000000
77	0.000000
78	0.000000
79	0.000000
80	0.000000
81	0.000000
82	0.000000
83	0.000000
84	0.000000
85	0.000000
86	0.000000
87	0.000000
88	0.000000
89	0.000000
90	0.000000
91	0.000000
92	0.000000
93	0.000000
94	0.000000
95	0.000000
96	0.000000
97	0.000000
98	0.000000
99	0.000000
100	0.000000
101	0.000000
102	0.000000
103	0.000000
104	0.000000
105	0.000000
106	0.000000
107	0.000000
108	0.000000
109	0.000000
110	0.000000
111	0.000000
112	0.000000
113	0.000000
114	0.000000
115	0.000000
116	0.000000
117	0.000000
118	0.000000
119	0.000000
120	0.000000
121	0.000000
122	0.000000
123	0.000000
124	0.000000
125	0.000000
126	0.000000
127	0.000000
128	0.000000
129	0.000000
130	0.000000
131	0.000000
132	0.000000
133	0.000000
134	0.000000
135	0.000000
136	0.000000
137	0.000000
138	0.000000
139	0.000000
140	0.000000
141	0.000000
142	0.000000
143	0.000000
144	0.000000
145	0.000000
146	0.000000
147	0.000000
148	0.000000
149	0.000000
150	0.000000
151	0.000000
152	0.000000
153	0.000000
154	0.000000
155	0.000000
156	0.000000
157	0.000000
158	0.000000
159	0.000000
160	0.000000
161	0.000000
162	0.000000
163	0.000000
164	0.000000
165	0.000000
166	0.000000
167	0.000000
168	0.000000
169	0.000000
170	0.000000
171	0.000000
172	0.000000
173	0.000000
174	0.000000
175	0.000000
176	0.000000
177	0.000000
178	0.000000
179	0.000000
180	0.000000
181	0.000000
182	0.000000
183	0.000000
184	0.000000
185	0.000000
186	0.000000
187	0.000000
188	0.000000
189	0.000000
190	0.000000
191	0.000000
192	0.000000
193	0.000000
194	0.000000
195	0.000000
196	0.000000
197	0.000000
198	0.000000
199	0.000000
200	0.000000
201	0.000000
202	0.000000
203	0.000000
204	0.000000
205	0.000000
206	0.000000
207	0.000000
208	0.000000
209	0.000000
210	0.000000
211	0.000000
212	0.000000
213	0.000000
214	0.000000
215	0.000000
216	0.000000
217	0.000000
218	0.000000
219	0.000000
220	0.000000
221	0.000000
222	0.000000
223	0.000000
224	0.000000
225	0.000000
226	0.000000
227	0.000000
228	0.000000
229	0.000000
230	0.000000
231	0.000000
232	0.000000
233	0.000000
234	0.000000
235	0.000000
236	0.000000
237	0.000000
238	0.000000
239	0.000000
240	0.000000
241	0.000000
242	0.000000
243	0.000000
244	0.000000
245	0.000000
246	0.000000
247	0.000000
248	0.000000
249	0.000000
250	0.000000
251	0.000000
252	0.000000
253	0.000000
254	0.000000
255	0.000000
256	0.000000
257	0.000000
258	0.000000
259	0.000000
260	0.000000
261	0.000000
262	0.000000
263	0.000000
264	0.000000
265	0.000000
266	0.000000
267	0.000000
268	0.000000
269	0.000000
270	0.000000
271	0.000000
272	0.000000
273	0.000000
274	0.000000
275	0.000000
276	0.000000
277	0.000000
278	0.000000
279	0.000000
280	0.000000
281	0.000000
282	0.000000
283	0.000000
284	0.000000
285	0.000000
286	0.000000
287	0.000000
288	0.000000
289	0.000000
290	0.000000
291	0.000000
292	0.0



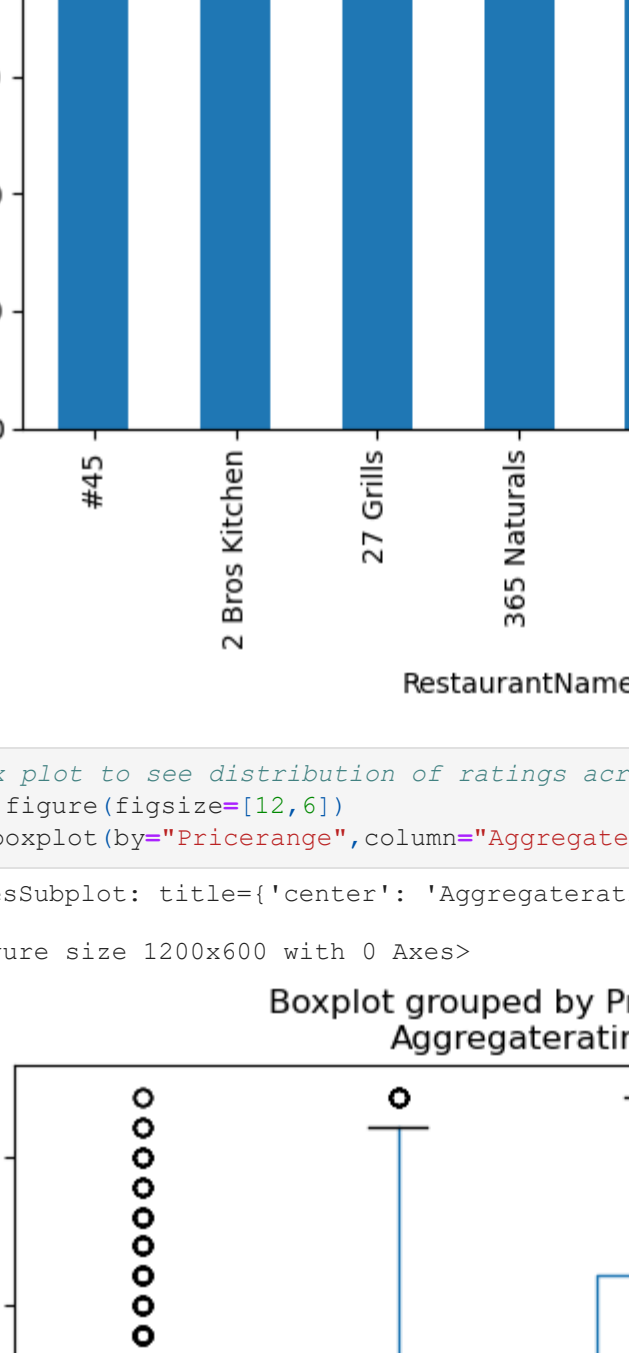
```
In [120]: #Which country makes up the most for somato restaurants
rest=rr['Country'].unique()
numbers=rr['Country'].value_counts()

plt.pie(numbers, labels = rest, startangle = 75, autopct='%1.0f%%')
plt.show()
```



```
In [39]: #Are delivery restaurants more or not?
rest=rr['HasOnlineDelivery'].unique()
numbers=rr['HasOnlineDelivery'].value_counts()

plt.pie(numbers, labels = rest, startangle = 75, autopct='%1.2f%%')
plt.show()
```



```
In [48]: #Are table booking restaurants more or not?
bookyes=len(rr[rr['HasTablebooking']=='Yes'])
bookno=len(rr[rr['HasTablebooking']=='No'])

ratio=bookyes/bookno

print("The ratio to restaurants that allow bookings vs dont allow",ratio)

#The ratio of restaurants that allow bookings vs dont allow is ~0.1381531853972799
```

```
In [41]: #Difference of votes btw delivery non delivery restaurants
bookyes= rr[rr['HasOnlineDelivery']=='Yes']
bookno=rr[rr['HasOnlineDelivery']=='No']
bookyesvote= bookyes['Votes'].sum()
booknovote= bookno['Votes'].sum()
diff=bookyesvote-booknovote
print("The difference of votes between delivery and non delivery restaurants is",diff)

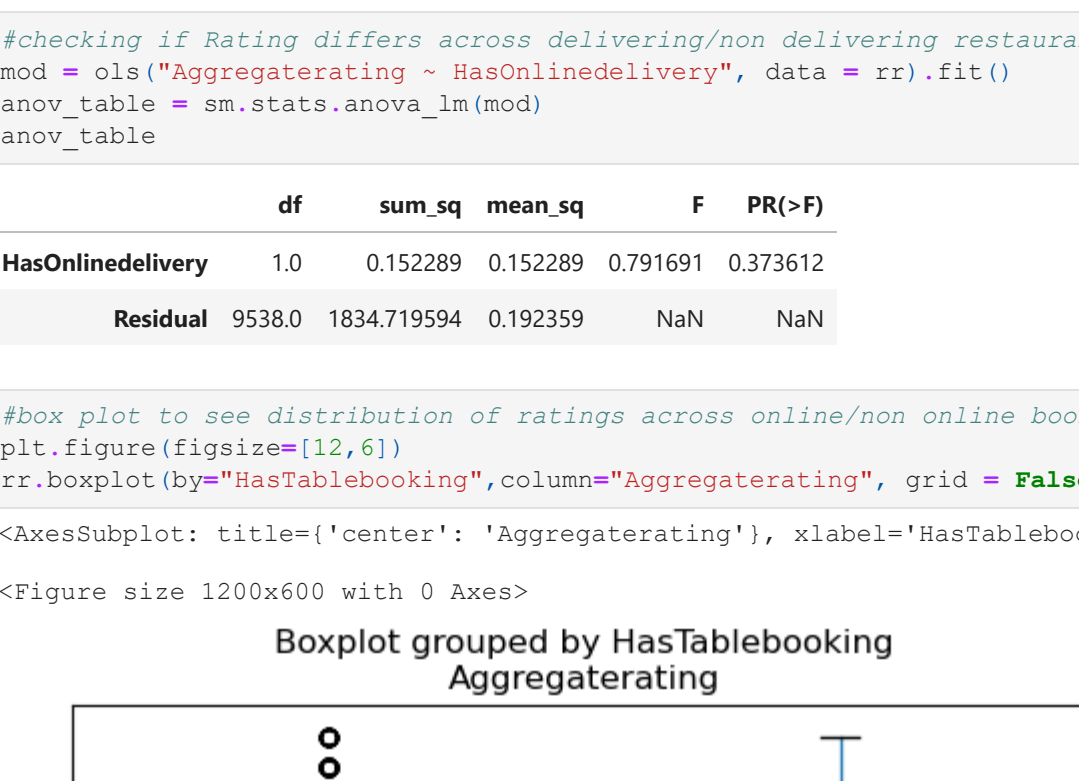
#The difference of votes between delivery and non delivery restaurants is ~27374.0
```

```
In [42]: #Cuisine count
print("The most number of cuisines offered by a restaurant is",rr['cuisinecount'].max(),
      "and the least is",rr['cuisinecount'].min())

#The most number of cuisines offered by a restaurant is 8.0 and the least is 0.0
```

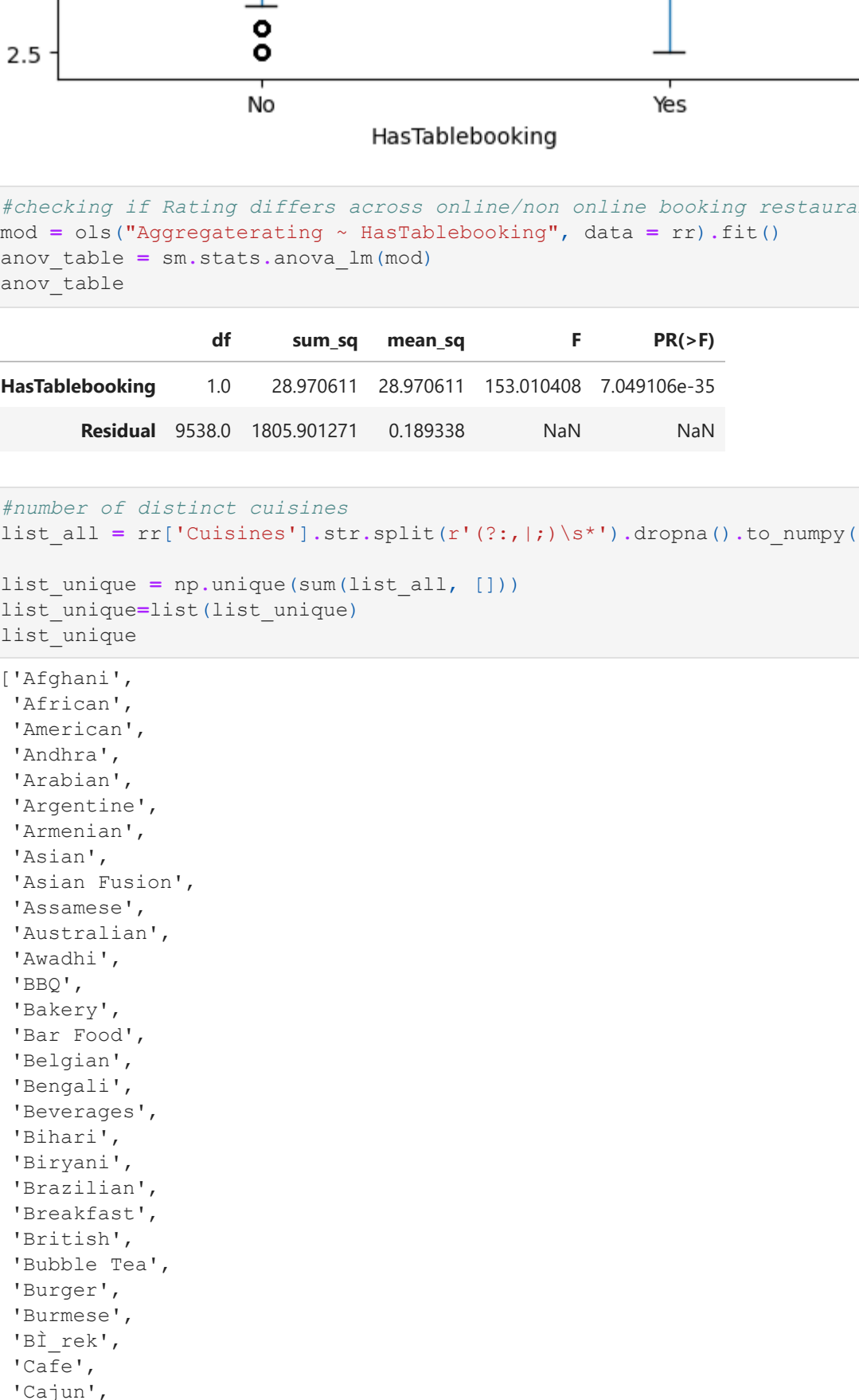
```
In [43]: #Average cost of meals across top restaurants
s_by_city= rr.groupby("RestaurantName")["AverageCostfortwo"].mean().nlargest(7)
s_by_city.plot(kind = "bar")

<AxesSubplot: xlabel='RestaurantName'>
```



```
In [54]: #Box plot to see distribution of ratings across price range
plt.figure(figsize=(12,6))
rr.boxplot(by="Pricerange",column="Aggregaterating", grid = False)

Out[54]: <Figure size 1200x600 with 0 Axes>
```



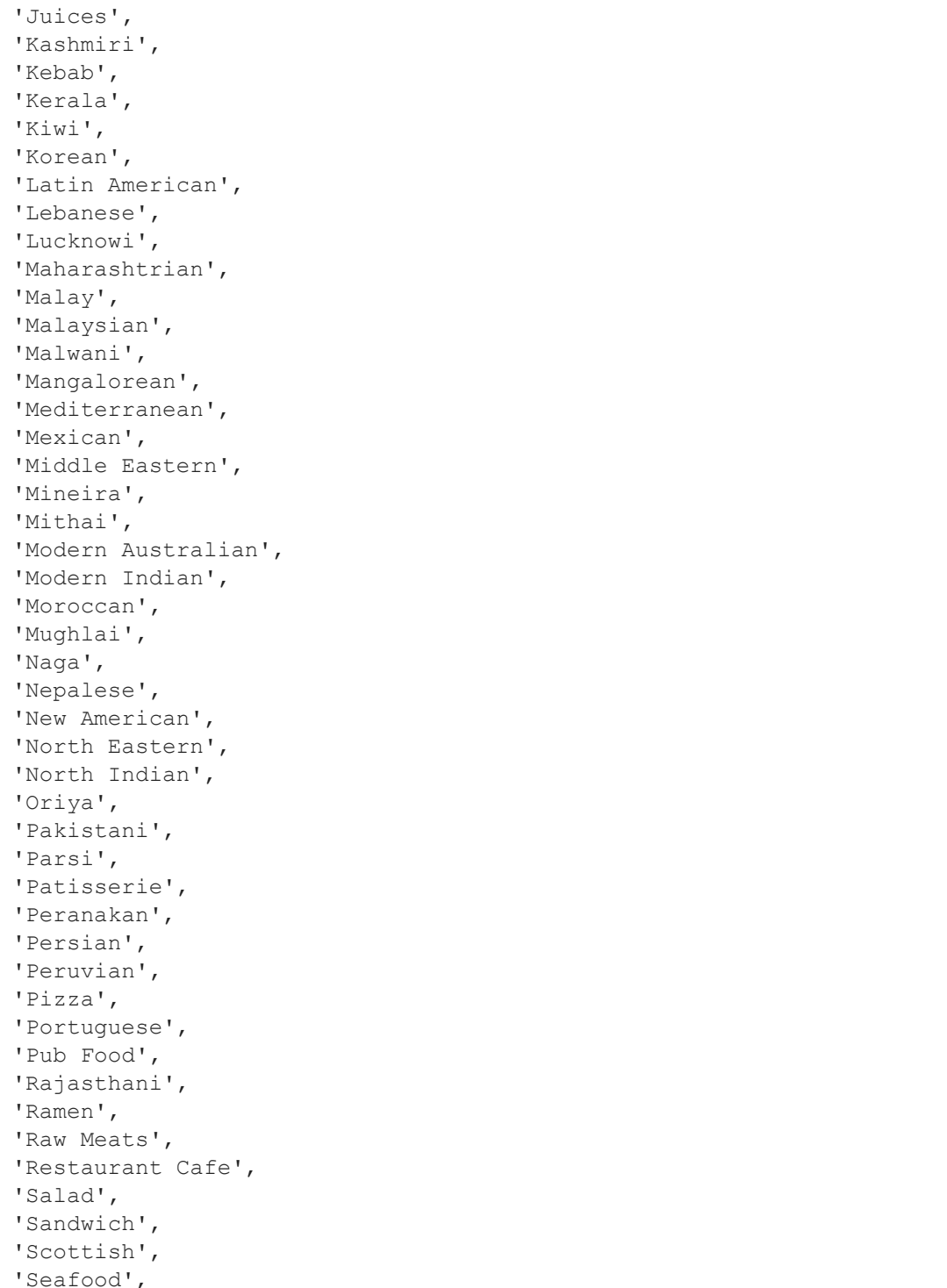
```
In [44]: #checking if Rating differs according to the Price range of the restaurant
mod = ols("Aggregaterating ~ Pricerange", data = rr).fit()
anov_table = sm.stats.anova_lm(mod)

Out[44]:
```

	df	sum_sq	mean_sq	F	PR(>F)
Pricerange	1.0	264.426240	264.426240	1605.975663	0.0
Residual	9538.0	1570.445642	0.164851	NaN	NaN

```
In [55]: #Box plot to see distribution of ratings across delivering/non delivering restaurants
plt.figure(figsize=(12,6))
rr.boxplot(by="HasOnlineDelivery",column="Aggregaterating", grid = False)

Out[55]: <Figure size 1200x600 with 0 Axes>
```



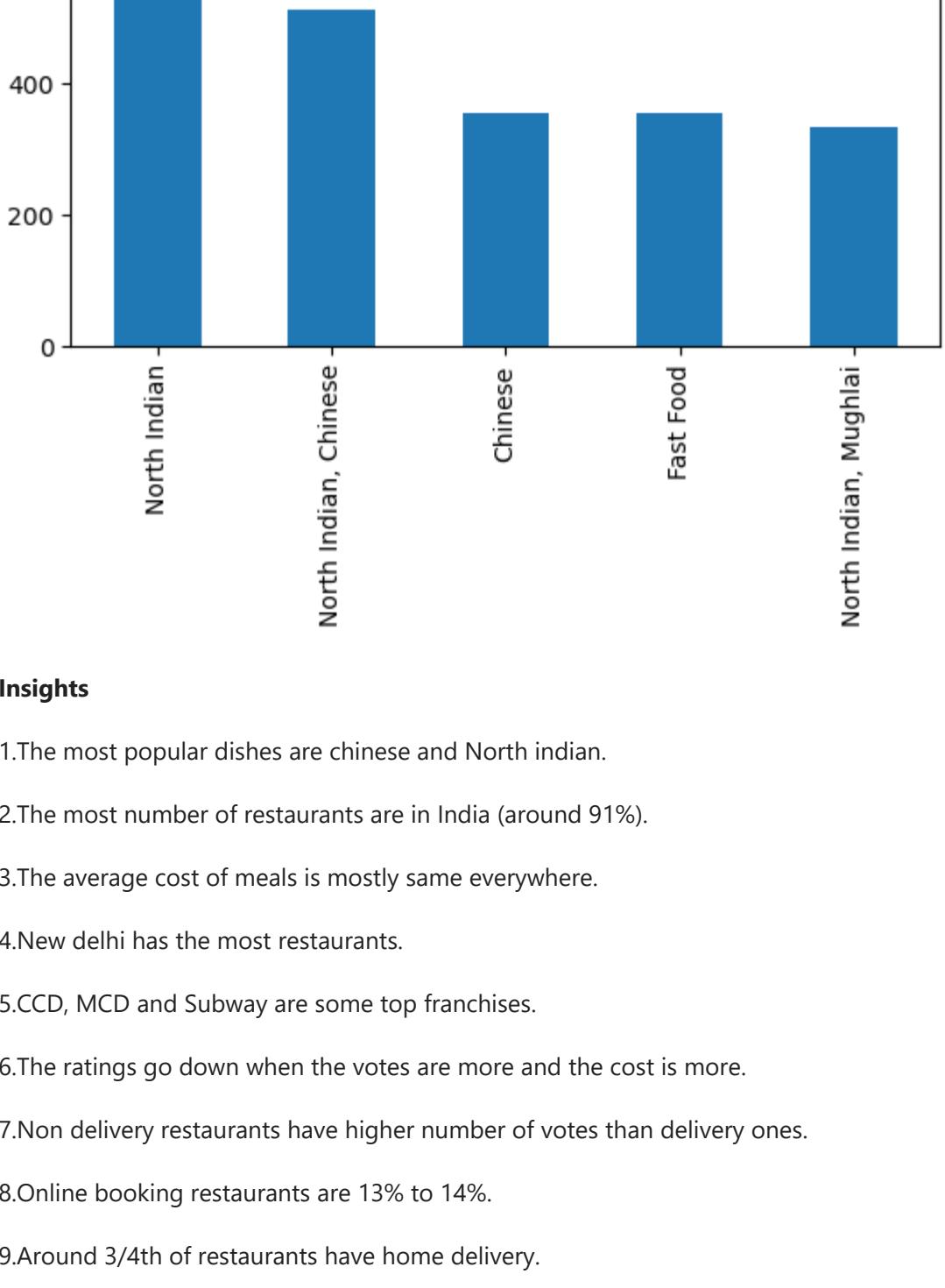
```
In [56]: #checking if Rating differs across delivering/non delivering restaurants
mod = ols("Aggregaterating ~ HasOnlineDelivery", data = rr).fit()
anov_table = sm.stats.anova_lm(mod)

Out[56]:
```

	df	sum_sq	mean_sq	F	PR(>F)
HasOnlineDelivery	1.0	0.152289	0.152289	0.791691	0.373612
Residual	9538.0	1834.719594	0.192359	NaN	NaN

```
In [57]: #Box plot to see distribution of ratings across online/non online booking restaurants
plt.boxplot(
    rr.boxplot(
        by="HasTablebooking",column="Aggregaterating", grid = False
    )
)

Out[57]: <Figure size 1200x600 with 0 Axes>
```



```
In [58]: #checking if Rating differs across online/non online booking restaurants
mod = ols("Aggregaterating ~ HasTablebooking", data = rr).fit()
anov_table = sm.stats.anova_lm(mod)

Out[58]:
```

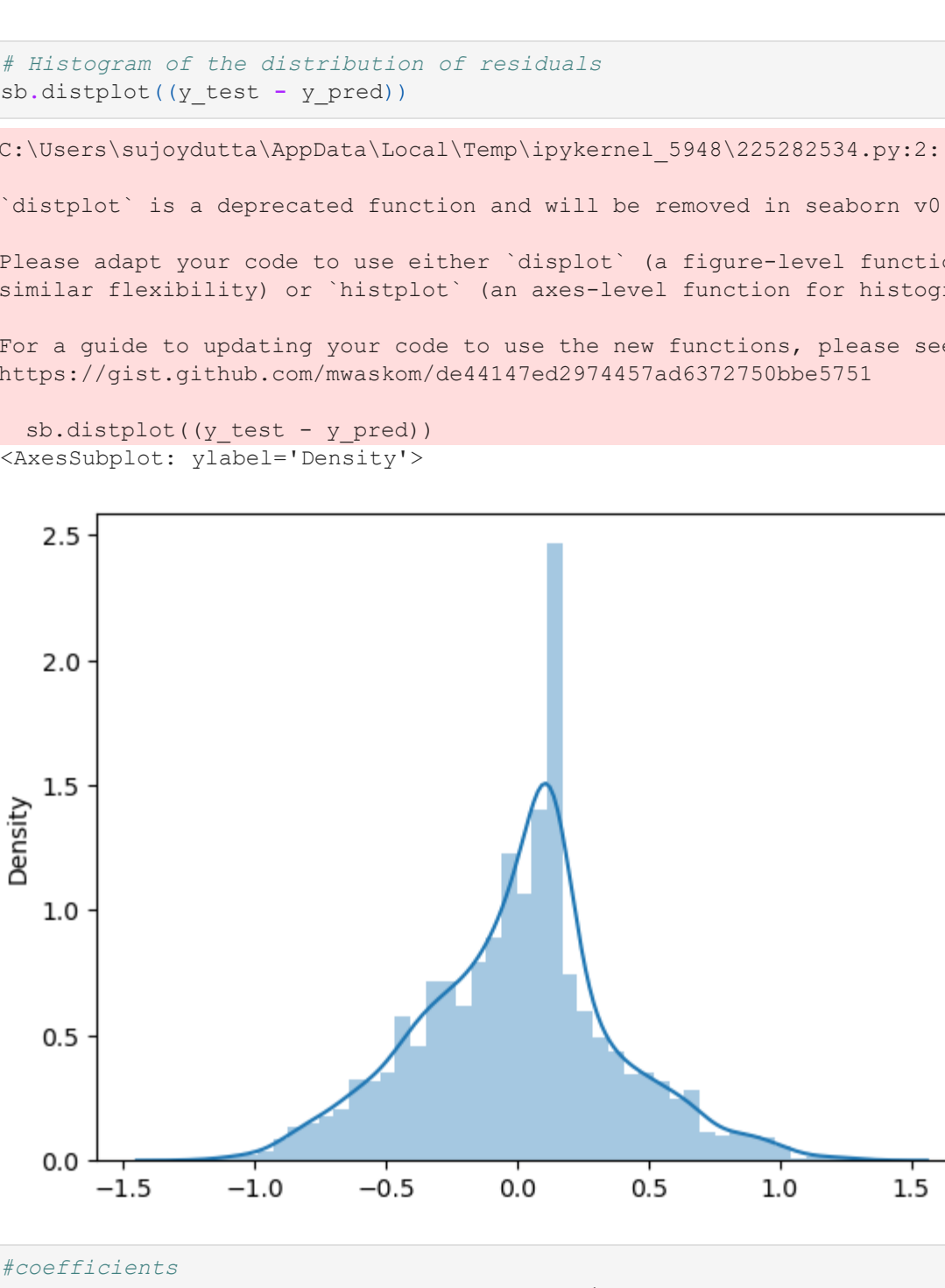
	df	sum_sq	mean_sq	F	PR(>F)
HasTablebooking	1.0	28.970611	28.970611	153.010408	7.049106e-35
Residual	9538.0	1805.901271	0.189338	NaN	NaN

```
In [90]: #Number of distinct cuisines
list_all = rr['Cuisines'].str.split("?",expand=True).dropna().to_numpy()
list_unique = np.unique(list_all, axis=0)
list_unique=list(list_unique)
list_unique
```

- 'Afghan',
- 'African',
- 'American',
- 'Andhra',
- 'Arabian',
- 'Argentine',
- 'Armenian',
- 'Asian',
- 'Asian Fusion',
- 'Awamese',
- 'Austrian',
- 'Awadhi',
- 'BBQ',
- 'Bakery',
- 'Bar Food',
- 'Belgian',
- 'Bengali',
- 'Beverages',
- 'Biryani',
- 'Brazilian',
- 'Breakfast',
- 'British',
- 'Bubble Tea',
- 'Burger',
- 'Burmese',
- 'Bake',
- 'Cajun',
- 'Canadian',
- 'Cantonese',
- 'Caribbean',
- 'Charcoal Grill',
- 'Chettinad',
- 'Chinese',
- 'Coffee and Tea',
- 'Contemporary',
- 'Continental',
- 'Cuban',
- 'Cuisine Varies',
- 'Curry',
- 'Deli',
- 'Desserts',
- 'Dish Sam',
- 'Diner',
- 'Drinks Only',
- 'Durban',
- 'DI_meat',
- 'Fajita',
- 'Fast Food',
- 'Filipino',
- 'Finger Food',
- 'Fish and Chips',
- 'French',
- 'Fusion',
- 'German',
- 'Goan',
- 'Gourmet Fast Food',
- 'Greek',
- 'Grill',
- 'Gujiati',
- 'Hawaiian',
- 'Healthy Food',
- 'Hyderabadi',
- 'Ice Cream',
- 'Indian',
- 'Indonesian',
- 'International',
- 'Iranian',
- 'Irish',
- 'Italian',
- 'Jigara',
- 'Japanese',
- 'Juices',
- 'Kashmiri',
- 'Kebab',
- 'Kerala',
- 'Kiwi',
- 'Korean',
- 'Latin American',
- 'Lebanese',
- 'Ladakh',
- 'Maharashtrian',
- 'Malay',
- 'Malaysian',
- 'Malian',
- 'Mangalorean',
- 'Mediterranean',
- 'Mexican',
- 'Middle Eastern',
- 'Mishai',
- 'Modern Australian',
- 'Modern Indian',
- 'Moroccan',
- 'Mughlai',
- 'Naga',
- 'Nepalese',
- 'New American',
- 'North Eastern',
- 'North Indian',
- 'Oriya',
- 'Pakistani',
- 'Parsi',
- 'Patisserie',
- 'Paranakan',
- 'Persian',
- 'Peruvian',
- 'Pizzeria',
- 'Portuguese',
- 'Pub Food',
- 'Rajasthani',
- 'Ramen',
- 'Raw Meats',
- 'Restaurant Cafe',
- 'Salad',
- 'Sandwich',
- 'Scottish',
- 'Seafood',
- 'Singaporean',
- 'Soul Food',
- 'South African',
- 'South American',
- 'South Indian',
- 'Southern',
- 'Southwestern',
- 'Spanish',
- 'Sri Lankan',
- 'Steak',
- 'Street Food',
- 'Sunda',
- 'Sushi',
- 'Szechuan',
- 'Tapa',
- 'Tea',
- 'Teryaki',
- 'Tex-Mex',
- 'Thai',
- 'Tibetan',
- 'Turkish',
- 'Turkish Pizza',
- 'Vegetarian',
- 'Vietnamese',
- 'Western',
- 'World Cuisine'

```
In [94]: #most popular cuisine type
from collections import Counter
from collections import Counter
from collections import Counter

Out[94]: <AxesSubplot: >
```



Insights

- The most popular dishes are chinese and north indian.
- The most number of restaurants are in India (around 91%).
- The average cost of meals is mostly same everywhere.
- New delhi has the most restaurants.
- CCD, MCD and Subway are some top franchises.
- The ratings go down when the votes are more and the cost is more.
- Non delivery restaurants have higher number of votes than delivery ones.
- Online booking restaurants are 13% to 14%.
- Around 3/4th of restaurants have home delivery.
- There are around 143 different cuisines.
- The ratings doesnt depend on home delivery options but on table booking options and price level.
- Restaurants can offer upto at most 8 types of cuisines.

Step 4: Model Building

After the data has been cleaned and formatted. Now its time to analyse and get insights. We will use Linear Regression to get the factors affecting ratings.

```
In [95]: rr.head()
```

	RestaurantName	CountryCode	City	Cuisines	cuisinecount	AverageCostfortwo	Currency	HasTablebooking	HasOnlineDelivery	Pricerange
0	Amber	1	New Delhi	North Indian, Chinese, Mughlai	3.0	300.0	Indian Rupees(INR)	Yes	Yes	
1	Bero's	1	New Delhi	Chinese, Thai	2.0	300.0	Indian Rupees(INR)	Yes	Yes	
2	Colone's Kababz	1	New Delhi	North Indian, Mughlai	2.0	300.0	Indian Rupees(INR)	Yes	No	
3	Oiva - The Italian Restaurant	1	New Delhi	Italian	1.0	300.0	Indian Rupees(INR)	Yes	Yes	
4	Drums of Heaven	1	New Delhi	Chinese, Seafood, Thai	3.0	300.0	Indian Rupees(INR)	Yes	Yes	

```
In [98]: #dropping not needed variables
rr1=rr.drop(['RestaurantName','City','Ratingtext'],axis=1)
rr1.head()
```

	CountryCode	Cuisines	cuisinecount	AverageCostfortwo	Currency	HasTablebooking	HasOnlineDelivery	Pricerange	Aggregaterating	Votes
--	-------------	----------	--------------	-------------------	----------	-----------------	-------------------	------------	-----------------	-------

```
In [96]: #creating target variable
target=rr1['Aggregaterating']
```

```
In [97]: #removing target variable from main dataframe
rr1=rr.drop(['Aggregaterating'], axis=1)
```

```
In [101]: #defining category variables
category_variables=['CountryCode', 'cuisinecount','Cuisines',
                    'HasTablebooking', 'HasOnlineDelivery', 'Pricerange', 'Country']
category_variables
```

```
In [101]: #CountryCode',
'cuisinecount',
'Cuisines',
'HasTablebooking',
'HasOnlineDelivery',
'Pricerange',
'Country']

In [105]: #initializing encoder
label_encoder=preproc.LabelEncoder()
```

```
In [106]: #encoding categorical variables
cuisenc=label_encoder.fit_transform(rr1['Cuisines'])
odenc=label_encoder.fit_transform(rr1['HasOnlineDelivery'])
ctenc=label_encoder.fit_transform(rr1['Country'])
```

```
In [107]: # Create linear regression object
regressor = LinearRegression()
```

```
In [110]: #setting values of X and y
X=np.d.DataFrame([cuisenc,
                  odenc,
                  ctenc,
                  rr1['CountryCode'],
                  rr1['cuisinecount'],
                  rr1['Pricerange']])
y= target.values
```

```
In [111]: # Train, test, split
X_train,X_test,y_train,y_test = train_test_split(X,y, test_size = .30, random_state= 1000)
```

```
In [112]: # Fit model to training data
regressor.fit(X_train,y_train)

Out[112]:
```

```
# Predict
# Predicting test set results
y_pred = regressor.predict(X_test)
```

```
Out[113]: array([3.79473488, 3.44930975, 3.57030208, ..., 3.28853989, 3.76064075,
3.2468455 ])
```

```
In [114]: # Calculated R Squared
print("R^2 =",metrics.explained_variance_score(y_test,y_pred))

R^2 = 0.257181190893376
```

```
In [115]: # Actual v predictions scatter
plt.scatter(y_test,y_pred)

Out[115]: <matplotlib.collections.PathCollection at 0x1860c433e20>
```



```
In [116]: # Histogram of the distribution of residuals
sb.distplot((y_test - y_pred))

Out[116]:
```



```
In [117]: #Coefficients
cof = pd.DataFrame(data = regressor.coef_, index = X.columns, columns = ['Coefficients'])
cof
```

	Coefficients
0	-0.000079
1	0.016249
2	-0.024005
3	-0.015809
4	0.003059
5	-0.022538
6	0.155563

Conclusion

The model is not very good as it can explain only 25% of the change in the dependent variable. The fit is not achieved. The variables are not enough and maybe another methodology needs to be used to develop this model.