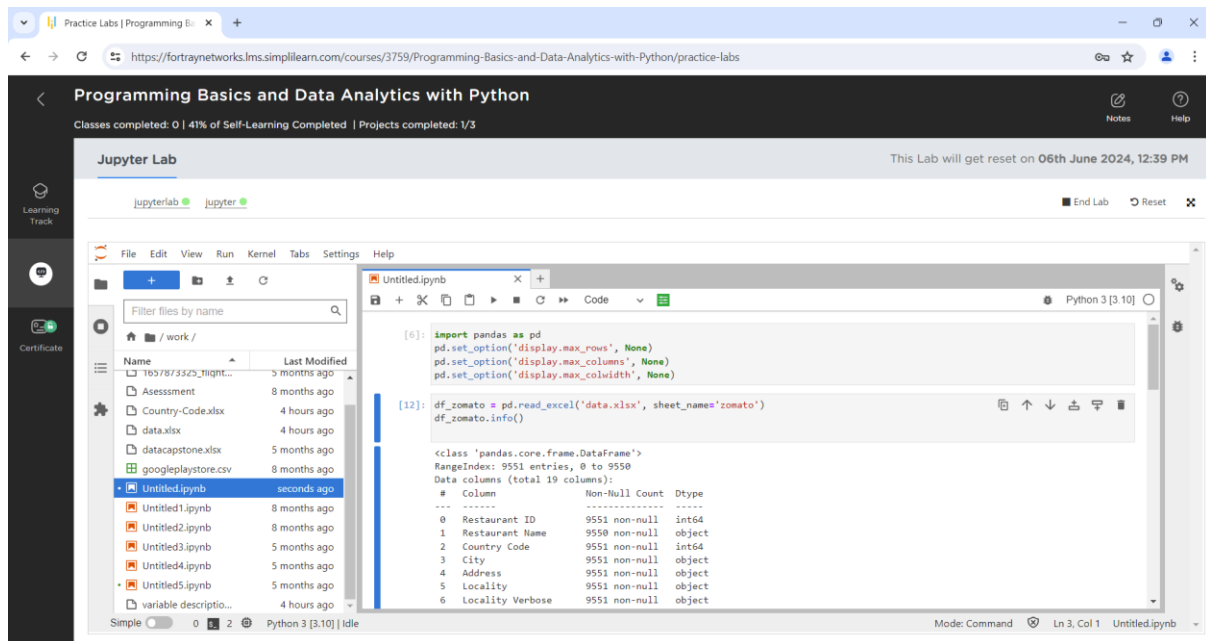Data Preliminary analysis:

Perform preliminary data inspection and report the findings as the structure of the data, missing values,

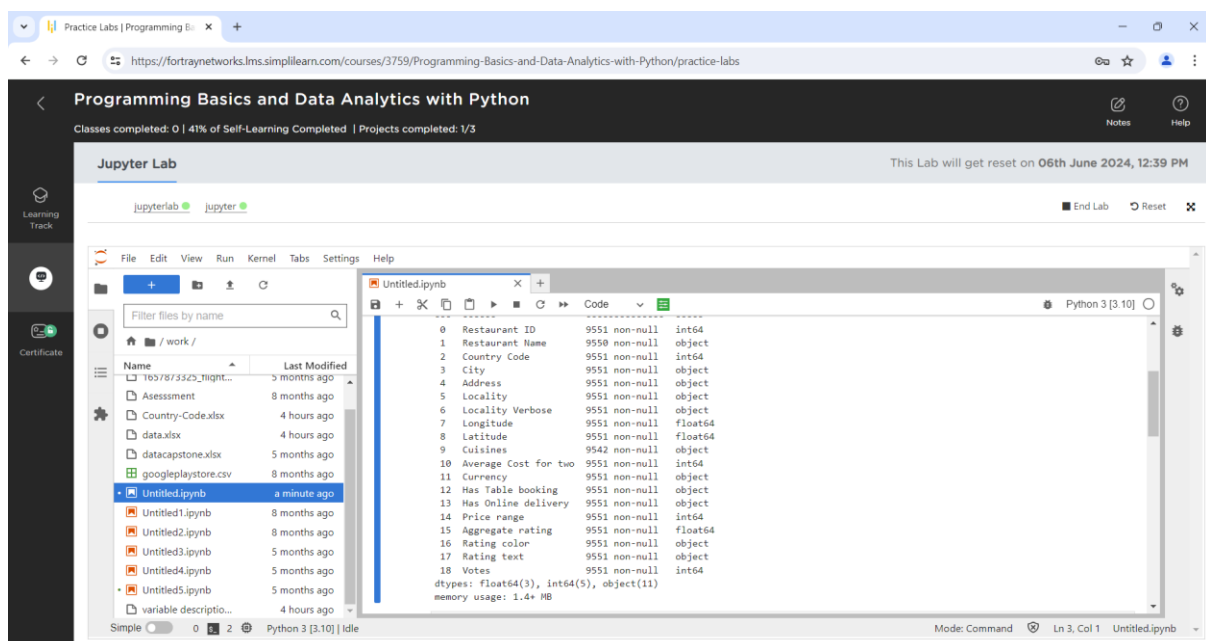duplicates cleaning variable names etc.

Import packages for data visualization:



Convert all column names to lowercase and use underscores to separate words.

Rename the "votes" column to "item_votes:



Check null values:

Looks like there are some null values in the columns restaurant_name and cuisines

Fill the null values values in the columns restaurant_name and cuisines with "Not Available"



# Check duplicate rows:

Explore the geographical distribution of the restaurants, finding out the cities with maximum / minimum number of restaurants.

Explore how ratings are distributed overall.

Explore the franchise with most national presence:

What is the ratio between restaurants that allow table booking vs that do not allow table booking?

What is the percentage of restaurants providing online delivery?



```python
[25]:  # visualize the ration of restaurants with online delivery
        online_delivery_counts = df_zomato.groupby('has_online_delivery')['restaurant_name'].nunique()

        # set plot size
        plt.figure(figsize=(4, 4))

        # plot pie chart
        plt.pie(online_delivery_counts, labels=online_delivery_counts.index, autopct='%1.1f%%', startangle=140, colors=sns.color_pale

        # add title
        plt.title('Distribution of online delivery')

        # show the plot
        plt.show()
```

Distribution of online delivery

Is there a difference in no. of votes for the restaurants that deliver and the restaurant that don't?

What are the top 10 cuisines served across cities?



What is the maximum and minimum no. of cuisines that a restaurant serves? Also, what is the relationship between No. of cuisines served and Ratings

**Programming Basics and Data Analytics with Python**

Classes completed: 0 | 41% of Self-Learning Completed | Projects completed: 1/3

Notes   Help

**Jupyter Lab**

This Lab will get reset on **06th June 2024, 12:39 PM**

jupyterlab ● jupyter ●

Used 4.7 of 50 hours in May, 2024    ■ End Lab    ↻ Reset

Learning Track

Certificate

File  Edit  View  Run  Kernel  Tabs  Settings  Help

Untitled.ipynb    ×    Untitled5.ipynb    ×    +

Code    Python 3 [3.10]

Filter files by name

🏠 ▪ / work /

Name: cuisines, dtype: int64

```
[88]: df_zomato_all_cuisines = df_zomato[["restaurant_id", "restaurant_name", "cuisines", "rating_text", "aggregate_rating"]].copy(
```

```
[89]: df_zomato_all_cuisines["len_cuisines"] = df_zomato_all_cuisines["cuisines"].str.len()
```

```
[90]: df_zomato_all_cuisines
```

| Name | Last Modified |
|---|---|
| lib | 5 months ago |
| 1657873325_flight... | 5 months ago |
| Asesssment | 8 months ago |
| Country-Code.xlsx | 9 hours ago |
| data.xlsx | 9 hours ago |
| datacapstone.xlsx | 5 months ago |
| googleplaystore.csv | 8 months ago |
| Untitled.ipynb | seconds ago |
| Untitled1.ipynb | 8 months ago |
| Untitled2.ipynb | 8 months ago |
| Untitled3.ipynb | 5 months ago |
| Untitled4.ipynb | 5 months ago |
| Untitled5.ipynb | 5 months ago |

[90]:

| | restaurant_id | restaurant_name | cuisines | rating_text | aggregate_rating | len_cuisines |
|---|---|---|---|---|---|---|
| 0 | 7402935 | Skye | [Italian, Continental] | Very Good | 4.1 | 2 |
| 1 | 7410290 | Satoo - Hotel Shangri-La | [Asian, Indonesian, Western] | Excellent | 4.6 | 3 |
| 2 | 7420899 | Sushi Masa | [Sushi, Japanese] | Excellent | 4.9 | 2 |
| 3 | 7421967 | 3 Wise Monkeys | [Japanese] | Very Good | 4.2 | 1 |
| 4 | 7422489 | Avec Moi Restaurant and Bar | [French, Western] | Very Good | 4.3 | 2 |
| 5 | 18352452 | Lucky Cat Coffee & Kitchen | [Cafe, Western] | Very Good | 4.3 | 2 |

Simple    0    1    Python 3 [3.10] | Idle    Mode: Command    Ln 1, Col 1    Untitled.ipynb

---

**Programming Basics and Data Analytics with Python**

Classes completed: 0 | 41% of Self-Learning Completed | Projects completed: 1/3

Notes   Help

**Jupyter Lab**

This Lab will get reset on **06th June 2024, 12:39 PM**

jupyterlab ● jupyter ●

Used 4.7 of 50 hours in May, 2024    ■ End Lab    ↻ Reset

Learning Track

Certificate

File  Edit  View  Run  Kernel  Tabs  Settings  Help

Untitled.ipynb    ×    Untitled5.ipynb    ×    +

Code    Python 3 [3.10]

Filter files by name

🏠 ▪ / work /

```python
[91]: # select the columns
correlation_data = df_zomato_all_cuisines[["len_cuisines", "aggregate_rating"]].copy()

# calculate the correlation matrix
correlation_matrix = correlation_data.corr()

# set plot size
plt.figure(figsize=(4, 4))

# create the heatmap
sns.heatmap(correlation_matrix, annot=True, cmap='Reds', fmt=".2f")

# add title
plt.title('Correlation Heatmap')

# show the plot
plt.show()
```
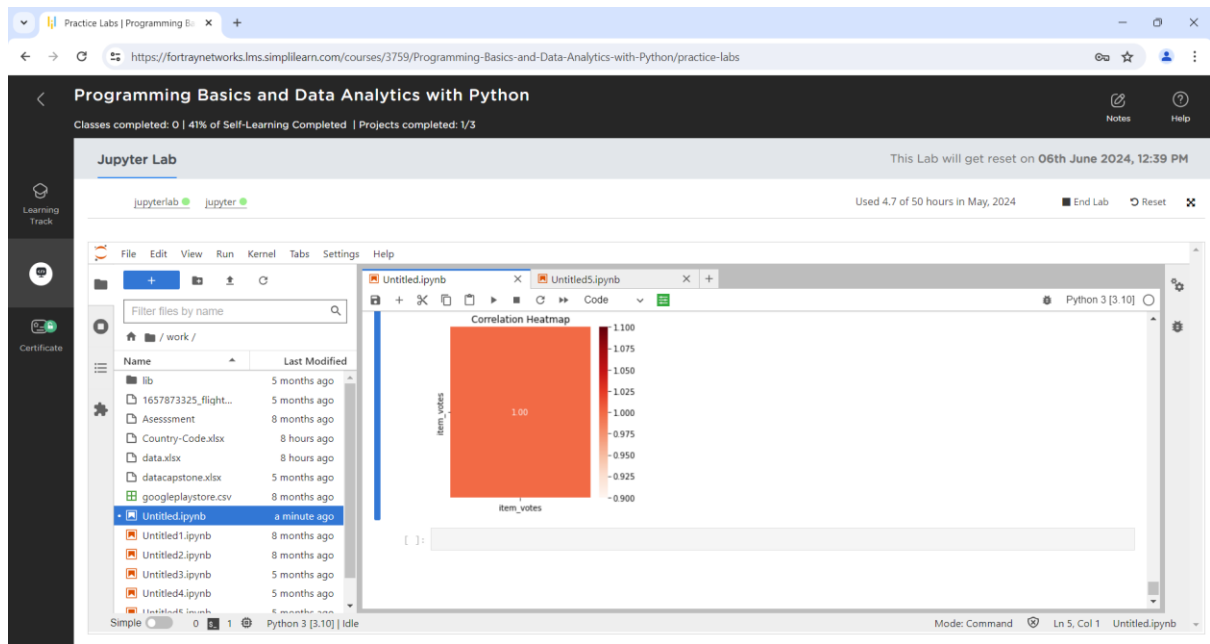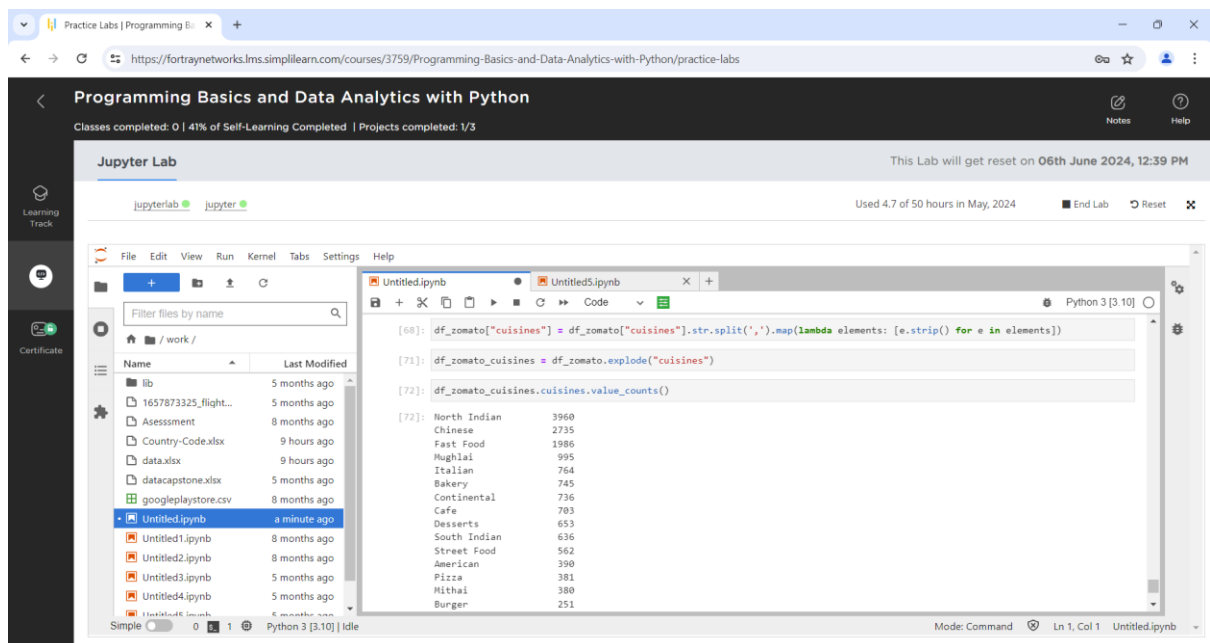
| Name | Last Modified |
|---|---|
| lib | 5 months ago |
| 1657873325_flight... | 5 months ago |
| Asesssment | 8 months ago |
| Country-Code.xlsx | 9 hours ago |
| data.xlsx | 9 hours ago |
| datacapstone.xlsx | 5 months ago |
| googleplaystore.csv | 8 months ago |
| Untitled.ipynb | in a few seconds |
| Untitled1.ipynb | 8 months ago |
| Untitled2.ipynb | 8 months ago |
| Untitled3.ipynb | 5 months ago |
| Untitled4.ipynb | 5 months ago |
| Untitled5.ipynb | 5 months ago |

Correlation Heatmap    1.0

Simple    0    1    Python 3 [3.10] | Idle    Saving completed    Mode: Command    Ln 1, Col 1    Untitled.ipynb

Discuss the cost vs the other variables.

It revealed that the median price for all items in India is 450 rupees, with the majority falling within the range of 300 to 700 rupees. However, several outliers are present, indicating dishes with significantly higher prices, with one reaching as high as 8000 rupees.

Explain the factors in the data that may have an effect on ratings e.g. No. of cuisines, cost, delivery option etc.

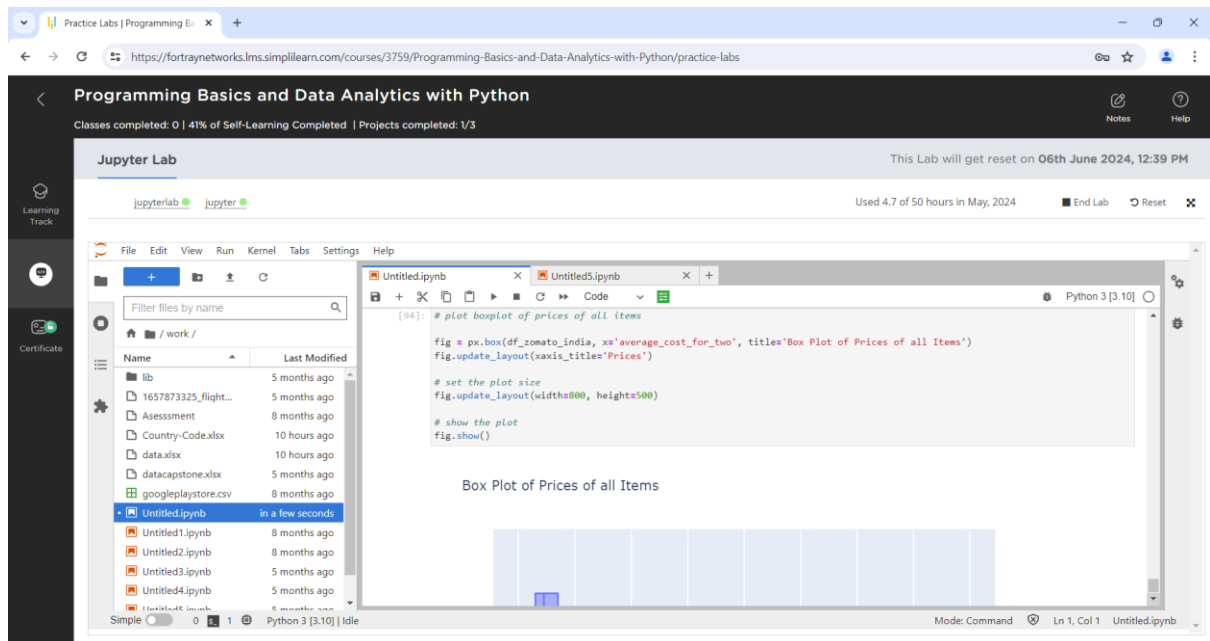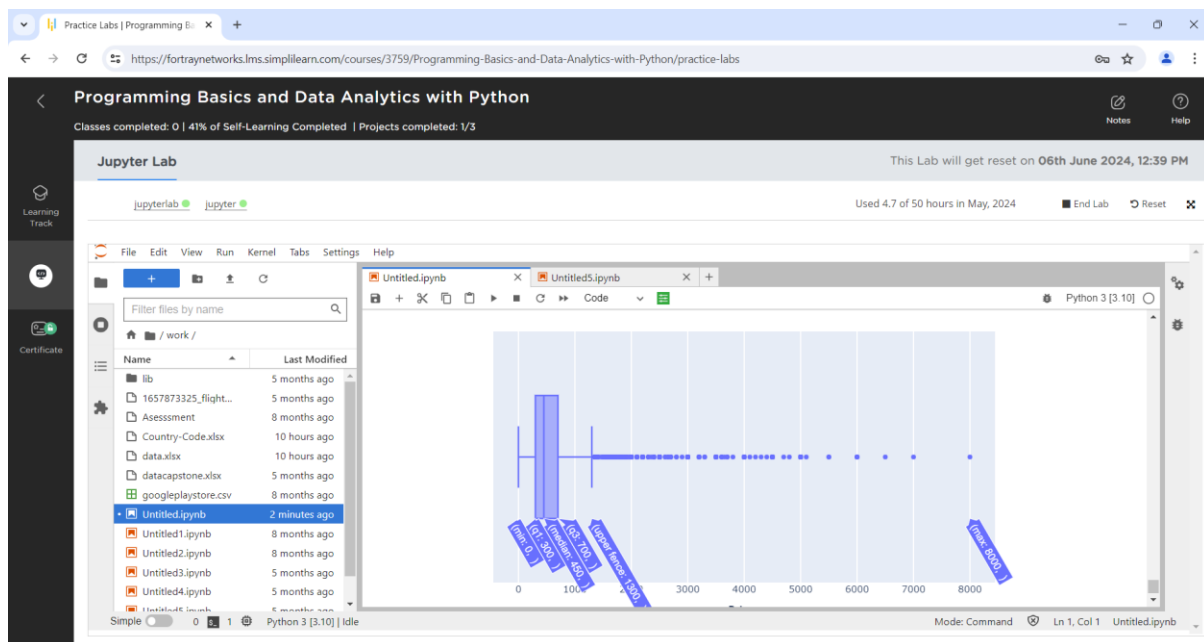Dining Type Analysis offered valuable insights into the customer's inclination towards dining or delivery services. Versatility in service, offering both dining and delivery options, proved to be a key factor in maximizing engagement.

The City Analysis revealed the best-suited locations for opening a restaurant, considering factors like engagement and competition. New Delhi surfaced as the busiest location, having maximum number of competitors.