# Activity_Course 2 Automatidata project lab

November 22, 2025

## 1 Automatidata project

**Course 2 - Get Started with Python**

Welcome to the Automatidata Project!

You have just started as a data professional in a fictional data consulting firm, Automatidata. Their client, the New York City Taxi and Limousine Commission (New York City TLC), has hired the Automatidata team for its reputation in helping their clients develop data-based solutions.

The team is still in the early stages of the project. Previously, you were asked to complete a project proposal by your supervisor, DeShawn Washington. You have received notice that your project proposal has been approved and that New York City TLC has given the Automatidata team access to their data. To get clear insights, New York TLC's data must be analyzed, key variables identified, and the dataset ensured it is ready for analysis.

A notebook was structured and prepared to help you in this project. Please complete the following questions.

## 2 Course 2 End-of-course project: Inspect and analyze data

In this activity, you will examine data provided and prepare it for analysis. This activity will help ensure the information is,

1. Ready to answer questions and yield insights

2. Ready for visualizations

3. Ready for future hypothesis testing and statistical methods

**The purpose** of this project is to investigate and understand the data provided.

**The goal** is to use a dataframe contructed within Python, perform a cursory inspection of the provided dataset, and inform team members of your findings.
*This activity has three parts:*

**Part 1:** Understand the situation * Prepare to understand and organize the provided taxi cab dataset and information.

**Part 2:** Understand the data

- Create a pandas dataframe for data learning, future exploratory data analysis (EDA), and statistical activities.

- Compile summary information about the data to inform next steps.

**Part 3:** Understand the variables

- Use insights from your examination of the summary data to guide deeper investigation into specific variables.

Follow the instructions and answer the following questions to complete the activity. Then, you will complete an Executive Summary using the questions listed on the PACE Strategy Document.

Be sure to complete this activity before moving on. The next course item will provide you with a completed exemplar to compare to your own work.

# 3 Identify data types and relevant variables using Python

# 4 PACE stages

Throughout these project notebooks, you'll see references to the problem-solving framework PACE. The following notebook components are labeled with the respective PACE stage: Plan, Analyze, Construct, and Execute.

## 4.1 PACE: Plan

Consider the questions in your PACE Strategy Document and those below to craft your response:

### 4.1.1 Task 1. Understand the situation

- How can you best prepare to understand and organize the provided taxi cab information?

==> ENTER YOUR RESPONSE HERE

## 4.2 PACE: Analyze

Consider the questions in your PACE Strategy Document to reflect on the Analyze stage.

### 4.2.1 Task 2a. Build dataframe

Create a pandas dataframe for data learning, and future exploratory data analysis (EDA) and statistical activities.

**Code the following,**

- import pandas as `pd`. pandas is used for buidling dataframes.

- import numpy as `np`. numpy is imported with pandas

- `df = pd.read_csv('Datasets\NYC taxi data.csv')`

**Note:** pair the data object name `df` with pandas functions to manipulate data, such as `df.groupby()`.

**Note:** As shown in this cell, the dataset has been automatically loaded in for you. You do not need to download the .csv file, or provide more code, in order to access the dataset and proceed with this lab. Please continue with this activity by completing the following instructions.

```python
[2]: #Import libraries and packages listed above
     ### YOUR CODE HERE ###
     import pandas as pd
     import numpy as np

     # Load dataset into dataframe
     df = pd.read_csv('2017_Yellow_Taxi_Trip_Data.csv')
     print("done")
```

```
done
```

### 4.2.2  Task 2b.  Understand the data - Inspect the data

View and inspect summary information about the dataframe by coding the following:

1. `df.head(10)`
2. `df.info()`
3. `df.describe()`

Consider the following two questions:

**Question 1:** When reviewing the `df.info()` output, what do you notice about the different variables? Are there any null values? Are all of the variables numeric? Does anything else stand out?

**Question 2:** When reviewing the `df.describe()` output, what do you notice about the distributions of each variable? Are there any questionable values?

==> ENTER YOUR RESPONSE TO QUESTIONS 1 & 2 HERE

```python
[3]: #==> ENTER YOUR CODE HERE
     df.head(10)
```

```
[3]:    Unnamed: 0  VendorID    tpep_pickup_datetime    tpep_dropoff_datetime  \
     0    24870114         2    03/25/2017 8:55:43 AM    03/25/2017 9:09:47 AM
     1    35634249         1    04/11/2017 2:53:28 PM    04/11/2017 3:19:58 PM
     2   106203690         1    12/15/2017 7:26:56 AM    12/15/2017 7:34:08 AM
     3    38942136         2    05/07/2017 1:17:59 PM    05/07/2017 1:48:14 PM
     4    30841670         2   04/15/2017 11:32:20 PM   04/15/2017 11:49:03 PM
     5    23345809         2    03/25/2017 8:34:11 PM    03/25/2017 8:42:11 PM
     6    37660487         2    05/03/2017 7:04:09 PM    05/03/2017 8:03:47 PM
     7    69059411         2    08/15/2017 5:41:06 PM    08/15/2017 6:03:05 PM
```

```
8    8433159           2   02/04/2017 4:17:07 PM    02/04/2017 4:29:14 PM
9   95294817           1   11/10/2017 3:20:29 PM    11/10/2017 3:40:55 PM

   passenger_count  trip_distance  RatecodeID store_and_fwd_flag  \
0                6           3.34           1                  N
1                1           1.80           1                  N
2                1           1.00           1                  N
3                1           3.70           1                  N
4                1           4.37           1                  N
5                6           2.30           1                  N
6                1          12.83           1                  N
7                1           2.98           1                  N
8                1           1.20           1                  N
9                1           1.60           1                  N

   PULocationID  DOLocationID  payment_type  fare_amount  extra  mta_tax  \
0           100           231             1         13.0    0.0      0.5
1           186            43             1         16.0    0.0      0.5
2           262           236             1          6.5    0.0      0.5
3           188            97             1         20.5    0.0      0.5
4             4           112             2         16.5    0.5      0.5
5           161           236             1          9.0    0.5      0.5
6            79           241             1         47.5    1.0      0.5
7           237           114             1         16.0    1.0      0.5
8           234           249             2          9.0    0.0      0.5
9           239           237             1         13.0    0.0      0.5

   tip_amount  tolls_amount  improvement_surcharge  total_amount
0        2.76           0.0                    0.3         16.56
1        4.00           0.0                    0.3         20.80
2        1.45           0.0                    0.3          8.75
3        6.39           0.0                    0.3         27.69
4        0.00           0.0                    0.3         17.80
5        2.06           0.0                    0.3         12.36
6        9.86           0.0                    0.3         59.16
7        1.78           0.0                    0.3         19.58
8        0.00           0.0                    0.3          9.80
9        2.75           0.0                    0.3         16.55
```

```python
#==> ENTER YOUR CODE HERE
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 22699 entries, 0 to 22698
Data columns (total 18 columns):
 #   Column                 Non-Null Count  Dtype
---  ------                 --------------  -----
```

```
 0   Unnamed: 0             22699 non-null  int64
 1   VendorID              22699 non-null  int64
 2   tpep_pickup_datetime  22699 non-null  object
 3   tpep_dropoff_datetime 22699 non-null  object
 4   passenger_count       22699 non-null  int64
 5   trip_distance         22699 non-null  float64
 6   RatecodeID            22699 non-null  int64
 7   store_and_fwd_flag    22699 non-null  object
 8   PULocationID          22699 non-null  int64
 9   DOLocationID          22699 non-null  int64
 10  payment_type          22699 non-null  int64
 11  fare_amount           22699 non-null  float64
 12  extra                 22699 non-null  float64
 13  mta_tax               22699 non-null  float64
 14  tip_amount            22699 non-null  float64
 15  tolls_amount          22699 non-null  float64
 16  improvement_surcharge 22699 non-null  float64
 17  total_amount          22699 non-null  float64
dtypes: float64(8), int64(7), object(3)
memory usage: 3.1+ MB
```

[5]:
```python
#==> ENTER YOUR CODE HERE
df.describe()
```

[5]:

|       | Unnamed: 0   | VendorID     | passenger_count | trip_distance \ |
|-------|--------------|--------------|-----------------|-----------------|
| count | 2.269900e+04 | 22699.000000 | 22699.000000    | 22699.000000    |
| mean  | 5.675849e+07 | 1.556236     | 1.642319        | 2.913313        |
| std   | 3.274493e+07 | 0.496838     | 1.285231        | 3.653171        |
| min   | 1.212700e+04 | 1.000000     | 0.000000        | 0.000000        |
| 25%   | 2.852056e+07 | 1.000000     | 1.000000        | 0.990000        |
| 50%   | 5.673150e+07 | 2.000000     | 1.000000        | 1.610000        |
| 75%   | 8.537452e+07 | 2.000000     | 2.000000        | 3.060000        |
| max   | 1.134863e+08 | 2.000000     | 6.000000        | 33.960000       |

|       | RatecodeID   | PULocationID | DOLocationID | payment_type | fare_amount \ |
|-------|--------------|--------------|--------------|--------------|---------------|
| count | 22699.000000 | 22699.000000 | 22699.000000 | 22699.000000 | 22699.000000  |
| mean  | 1.043394     | 162.412353   | 161.527997   | 1.336887     | 13.026629     |
| std   | 0.708391     | 66.633373    | 70.139691    | 0.496211     | 13.243791     |
| min   | 1.000000     | 1.000000     | 1.000000     | 1.000000     | -120.000000   |
| 25%   | 1.000000     | 114.000000   | 112.000000   | 1.000000     | 6.500000      |
| 50%   | 1.000000     | 162.000000   | 162.000000   | 1.000000     | 9.500000      |
| 75%   | 1.000000     | 233.000000   | 233.000000   | 2.000000     | 14.500000     |
| max   | 99.000000    | 265.000000   | 265.000000   | 4.000000     | 999.990000    |

|       | extra        | mta_tax      | tip_amount   | tolls_amount \ |
|-------|--------------|--------------|--------------|----------------|
| count | 22699.000000 | 22699.000000 | 22699.000000 | 22699.000000   |
| mean  | 0.333275     | 0.497445     | 1.835781     | 0.312542       |
```
```

```
std          0.463097        0.039465        2.800626        1.399212
min         -1.000000       -0.500000        0.000000        0.000000
25%          0.000000        0.500000        0.000000        0.000000
50%          0.000000        0.500000        1.350000        0.000000
75%          0.500000        0.500000        2.450000        0.000000
max          4.500000        0.500000      200.000000       19.100000


       improvement_surcharge   total_amount
count           22699.000000   22699.000000
mean                0.299551      16.310502
std                 0.015673      16.097295
min                -0.300000    -120.300000
25%                 0.300000       8.750000
50%                 0.300000      11.800000
75%                 0.300000      17.800000
max                 0.300000    1200.290000
```

### 4.2.3   Task 2c. Understand the data - Investigate the variables

Sort and interpret the data table for two variables:`trip_distance` and `total_amount`.

**Answer the following three questions:**

**Question 1:** Sort your first variable (`trip_distance`) from maximum to minimum value, do the values seem normal?

**Question 2:** Sort by your second variable (`total_amount`), are any values unusual?

**Question 3:** Are the resulting rows similar for both sorts? Why or why not?

==> ENTER YOUR RESPONSES TO QUESTION 1-3 HERE

```
[10]: # ==> ENTER YOUR CODE HERE

      df.sort_values(by='trip_distance', ascending = False)

      # Sort the data by trip distance from maximum to minimum value
```

```
[10]:        Unnamed: 0   VendorID    tpep_pickup_datetime   tpep_dropoff_datetime  \
      9280     51810714          2   06/18/2017 11:33:25 PM   06/19/2017 12:12:38 AM
      13861    40523668          2    05/19/2017 8:20:21 AM    05/19/2017 9:20:30 AM
      6064     49894023          2   06/13/2017 12:30:22 PM    06/13/2017 1:37:51 PM
      10291    76319330          2   09/11/2017 11:41:04 AM   09/11/2017 12:18:58 PM
      29       94052446          2    11/06/2017 8:30:50 PM   11/07/2017 12:00:00 AM
      ...           ...        ...                      ...                      ...
      2440     63574825          1   07/26/2017 10:26:58 PM   07/26/2017 10:26:58 PM
      15916    47368116          1    06/29/2017 7:30:30 PM    06/29/2017 7:43:29 PM
      1350     91619825          2   10/30/2017 8:20:29 AM    10/30/2017 8:20:38 AM
      246      78660848          1    09/18/2017 8:50:53 PM    09/18/2017 8:51:03 PM
```

```
17788     58079289            1   07/08/2017 12:54:02 AM   07/08/2017 12:55:03 AM
```

|       | passenger_count | trip_distance | RatecodeID | store_and_fwd_flag \ |
|-------|-----------------|---------------|------------|----------------------|
| 9280  | 2               | 33.96         | 5          | N                    |
| 13861 | 1               | 33.92         | 5          | N                    |
| 6064  | 1               | 32.72         | 3          | N                    |
| 10291 | 1               | 31.95         | 4          | N                    |
| 29    | 1               | 30.83         | 1          | N                    |
| …     | …               | …             | …          | …                    |
| 2440  | 1               | 0.00          | 1          | N                    |
| 15916 | 1               | 0.00          | 1          | N                    |
| 1350  | 1               | 0.00          | 1          | N                    |
| 246   | 1               | 0.00          | 1          | N                    |
| 17788 | 2               | 0.00          | 1          | N                    |

|       | PULocationID | DOLocationID | payment_type | fare_amount | extra | mta_tax \ |
|-------|--------------|--------------|--------------|-------------|-------|-----------|
| 9280  | 132          | 265          | 2            | 150.00      | 0.0   | 0.0       |
| 13861 | 229          | 265          | 1            | 200.01      | 0.0   | 0.5       |
| 6064  | 138          | 1            | 1            | 107.00      | 0.0   | 0.0       |
| 10291 | 138          | 265          | 2            | 131.00      | 0.0   | 0.5       |
| 29    | 132          | 23           | 1            | 80.00       | 0.5   | 0.5       |
| …     | …            | …            | …            | …           | …     | …         |
| 2440  | 162          | 264          | 2            | 5.50        | 0.5   | 0.5       |
| 15916 | 79           | 148          | 3            | 8.50        | 1.0   | 0.5       |
| 1350  | 193          | 193          | 1            | 2.50        | 0.0   | 0.5       |
| 246   | 145          | 145          | 2            | 2.50        | 0.5   | 0.5       |
| 17788 | 158          | 158          | 3            | 2.50        | 0.5   | 0.5       |

|       | tip_amount | tolls_amount | improvement_surcharge | total_amount |
|-------|------------|--------------|-----------------------|--------------|
| 9280  | 0.00       | 0.00         | 0.3                   | 150.30       |
| 13861 | 51.64      | 5.76         | 0.3                   | 258.21       |
| 6064  | 55.50      | 16.26        | 0.3                   | 179.06       |
| 10291 | 0.00       | 0.00         | 0.3                   | 131.80       |
| 29    | 18.56      | 11.52        | 0.3                   | 111.38       |
| …     | …          | …            | …                     | …            |
| 2440  | 0.00       | 0.00         | 0.3                   | 6.80         |
| 15916 | 0.00       | 0.00         | 0.3                   | 10.30        |
| 1350  | 0.66       | 0.00         | 0.3                   | 3.96         |
| 246   | 0.00       | 0.00         | 0.3                   | 3.80         |
| 17788 | 0.00       | 0.00         | 0.3                   | 3.80         |

```
[22699 rows x 18 columns]
```

```
[11]: #==> ENTER YOUR CODE HERE

      df.sort_values(by='total_amount', ascending = False)
      df.head(20)
```

```
# Sort the data by total amount and print the top 20 values
```

[11]:
```
      Unnamed: 0  VendorID   tpep_pickup_datetime   tpep_dropoff_datetime  \
0       24870114         2   03/25/2017 8:55:43 AM   03/25/2017 9:09:47 AM
1       35634249         1   04/11/2017 2:53:28 PM   04/11/2017 3:19:58 PM
2      106203690         1   12/15/2017 7:26:56 AM   12/15/2017 7:34:08 AM
3       38942136         2   05/07/2017 1:17:59 PM   05/07/2017 1:48:14 PM
4       30841670         2  04/15/2017 11:32:20 PM  04/15/2017 11:49:03 PM
5       23345809         2   03/25/2017 8:34:11 PM   03/25/2017 8:42:11 PM
6       37660487         2   05/03/2017 7:04:09 PM   05/03/2017 8:03:47 PM
7       69059411         2   08/15/2017 5:41:06 PM   08/15/2017 6:03:05 PM
8        8433159         2   02/04/2017 4:17:07 PM   02/04/2017 4:29:14 PM
9       95294817         1   11/10/2017 3:20:29 PM   11/10/2017 3:40:55 PM
10      18017909         2  03/04/2017 11:58:00 AM  03/04/2017 12:13:12 PM
11      18600059         2   03/05/2017 7:15:30 PM   03/05/2017 7:52:18 PM
12      46782248         1   06/09/2017 7:00:26 PM   06/09/2017 7:20:11 PM
13      94113247         2  11/06/2017 11:35:05 PM  11/06/2017 11:42:57 PM
14      14168279         1   02/22/2017 3:18:31 PM   02/22/2017 3:42:50 PM
15      47444401         2   06/02/2017 6:41:39 AM   06/02/2017 6:57:47 AM
16      69088676         1   08/15/2017 7:48:08 PM   08/15/2017 8:00:37 PM
17      58691513         2   07/10/2017 1:36:31 PM   07/10/2017 1:48:43 PM
18      35388828         2   04/10/2017 6:12:58 PM   04/10/2017 6:17:39 PM
19      18383214         2   03/05/2017 4:01:07 AM   03/05/2017 4:14:11 AM

    passenger_count  trip_distance  RatecodeID store_and_fwd_flag  \
0                 6           3.34           1                  N
1                 1           1.80           1                  N
2                 1           1.00           1                  N
3                 1           3.70           1                  N
4                 1           4.37           1                  N
5                 6           2.30           1                  N
6                 1          12.83           1                  N
7                 1           2.98           1                  N
8                 1           1.20           1                  N
9                 1           1.60           1                  N
10                1           1.77           1                  N
11                2          18.90           2                  N
12                1           3.00           1                  N
13                1           2.39           1                  N
14                1           3.30           1                  N
15                1           5.93           1                  N
16                1           3.60           1                  N
17                2           1.71           1                  N
18                2           0.63           1                  N
19                2           2.77           1                  N
```

```
     PULocationID  DOLocationID  payment_type  fare_amount  extra  mta_tax  \
0             100           231             1         13.0    0.0      0.5
1             186            43             1         16.0    0.0      0.5
2             262           236             1          6.5    0.0      0.5
3             188            97             1         20.5    0.0      0.5
4               4           112             2         16.5    0.5      0.5
5             161           236             1          9.0    0.5      0.5
6              79           241             1         47.5    1.0      0.5
7             237           114             1         16.0    1.0      0.5
8             234           249             2          9.0    0.0      0.5
9             239           237             1         13.0    0.0      0.5
10            162           142             1         11.5    0.0      0.5
11            236           132             1         52.0    0.0      0.5
12             13           148             1         15.0    1.0      0.5
13            209            25             1          9.5    0.5      0.5
14            238           161             1         17.5    0.0      0.5
15            239           231             1         19.0    0.0      0.5
16            163            41             1         12.5    1.0      0.5
17            142           100             1          9.5    0.0      0.5
18            263           262             2          5.0    1.0      0.5
19             79            68             1         11.5    0.5      0.5

    tip_amount  tolls_amount  improvement_surcharge  total_amount
0         2.76          0.00                    0.3         16.56
1         4.00          0.00                    0.3         20.80
2         1.45          0.00                    0.3          8.75
3         6.39          0.00                    0.3         27.69
4         0.00          0.00                    0.3         17.80
5         2.06          0.00                    0.3         12.36
6         9.86          0.00                    0.3         59.16
7         1.78          0.00                    0.3         19.58
8         0.00          0.00                    0.3          9.80
9         2.75          0.00                    0.3         16.55
10        2.46          0.00                    0.3         14.76
11       14.58          5.54                    0.3         72.92
12        3.35          0.00                    0.3         20.15
13        2.16          0.00                    0.3         12.96
14        4.55          0.00                    0.3         22.85
15        3.00          0.00                    0.3         22.80
16        2.85          0.00                    0.3         17.15
17        0.00          0.00                    0.3         10.30
18        0.00          0.00                    0.3          6.80
19        3.20          0.00                    0.3         16.00
```

[12]:
```python
#==> ENTER YOUR CODE HERE
df.tail(n=20)
# Sort the data by total amount and print the bottom 20 values
```

```
[12]:         Unnamed: 0  VendorID    tpep_pickup_datetime   tpep_dropoff_datetime  \
       22679      627525         2   01/03/2017 7:17:55 AM   01/03/2017 7:24:09 AM
       22680    39375701         1   05/16/2017 6:35:44 AM   05/16/2017 6:35:51 AM
       22681    56013080         2   06/09/2017 6:24:49 PM   06/09/2017 6:36:15 PM
       22682    76140746         2   09/10/2017 4:55:27 PM   09/10/2017 5:09:43 PM
       22683    65803597         2   08/03/2017 5:30:04 PM   08/03/2017 5:41:52 PM
       22684    65789688         1   08/03/2017 4:36:32 PM   08/03/2017 4:46:23 PM
       22685    57495742         2  07/05/2017 10:42:46 PM  07/05/2017 10:49:29 PM
       22686     9039930         2   02/08/2017 6:13:26 PM   02/08/2017 7:34:11 PM
       22687    10648695         2   02/13/2017 3:09:25 PM   02/13/2017 3:15:37 PM
       22688    66450599         2   08/05/2017 9:23:29 PM   08/05/2017 9:26:11 PM
       22689    19137636         2  03/07/2017 12:25:52 PM  03/07/2017 12:39:40 PM
       22690    79394312         2   09/21/2017 1:44:42 PM   09/21/2017 1:52:06 PM
       22691   112546155        2   01/06/2017 1:50:14 AM   01/06/2017 1:56:47 AM
       22692    60425673         1   07/16/2017 3:22:51 AM   07/16/2017 3:40:52 AM
       22693    67858616         2  08/10/2017 10:20:04 PM  08/10/2017 10:29:31 PM
       22694    14873857         2   02/24/2017 5:37:23 PM   02/24/2017 5:40:39 PM
       22695    66632549         2   08/06/2017 4:43:59 PM   08/06/2017 5:24:47 PM
       22696    74239933         2   09/04/2017 2:54:14 PM   09/04/2017 2:58:22 PM
       22697    60217333         2  07/15/2017 12:56:30 PM   07/15/2017 1:08:26 PM
       22698    17208911         1   03/02/2017 1:02:49 PM   03/02/2017 1:16:09 PM

              passenger_count  trip_distance  RatecodeID store_and_fwd_flag  \
       22679                1           1.03           1                  N
       22680                1           1.30           1                  N
       22681                1           1.79           1                  N
       22682                3           2.16           1                  N
       22683                1           1.17           1                  N
       22684                2           1.20           1                  N
       22685                1           1.01           1                  N
       22686                5          10.64           1                  N
       22687                1           0.59           1                  N
       22688                3           0.44           1                  N
       22689                1           1.96           1                  N
       22690                1           0.89           1                  N
       22691                1           2.12           1                  N
       22692                1           5.70           1                  N
       22693                1           0.89           1                  N
       22694                3           0.61           1                  N
       22695                1          16.71           2                  N
       22696                1           0.42           1                  N
       22697                1           2.36           1                  N
       22698                1           2.10           1                  N

              PULocationID  DOLocationID  payment_type  fare_amount  extra  mta_tax  \
       22679            68            48             1          6.0    0.0      0.5
       22680           230           230             3          2.5    0.0      0.5
```

|       |     |     |   |      |     |     |
|-------|-----|-----|---|------|-----|-----|
| 22681 | 234 | 144 | 1 | 9.5  | 1.0 | 0.5 |
| 22682 | 142 | 68  | 1 | 11.0 | 0.0 | 0.5 |
| 22683 | 107 | 170 | 1 | 8.5  | 1.0 | 0.5 |
| 22684 | 68  | 50  | 1 | 8.0  | 1.0 | 0.5 |
| 22685 | 144 | 79  | 1 | 6.5  | 0.5 | 0.5 |
| 22686 | 170 | 70  | 1 | 52.0 | 1.0 | 0.5 |
| 22687 | 141 | 237 | 1 | 5.5  | 0.0 | 0.5 |
| 22688 | 230 | 163 | 2 | 4.0  | 0.5 | 0.5 |
| 22689 | 113 | 13  | 1 | 11.0 | 0.0 | 0.5 |
| 22690 | 43  | 142 | 1 | 7.0  | 0.0 | 0.5 |
| 22691 | 170 | 79  | 1 | 8.0  | 0.5 | 0.5 |
| 22692 | 249 | 17  | 1 | 19.0 | 0.5 | 0.5 |
| 22693 | 229 | 170 | 1 | 7.5  | 0.5 | 0.5 |
| 22694 | 48  | 186 | 2 | 4.0  | 1.0 | 0.5 |
| 22695 | 132 | 164 | 1 | 52.0 | 0.0 | 0.5 |
| 22696 | 107 | 234 | 2 | 4.5  | 0.0 | 0.5 |
| 22697 | 68  | 144 | 1 | 10.5 | 0.0 | 0.5 |
| 22698 | 239 | 236 | 1 | 11.0 | 0.0 | 0.5 |

|       | tip_amount | tolls_amount | improvement_surcharge | total_amount |
|-------|-----------|--------------|-----------------------|--------------|
| 22679 | 2.04      | 0.00         | 0.3                   | 8.84         |
| 22680 | 0.00      | 0.00         | 0.3                   | 3.30         |
| 22681 | 1.00      | 0.00         | 0.3                   | 12.30        |
| 22682 | 2.36      | 0.00         | 0.3                   | 14.16        |
| 22683 | 2.06      | 0.00         | 0.3                   | 12.36        |
| 22684 | 1.95      | 0.00         | 0.3                   | 11.75        |
| 22685 | 1.56      | 0.00         | 0.3                   | 9.36         |
| 22686 | 14.84     | 5.54         | 0.3                   | 74.18        |
| 22687 | 1.26      | 0.00         | 0.3                   | 7.56         |
| 22688 | 0.00      | 0.00         | 0.3                   | 5.30         |
| 22689 | 2.36      | 0.00         | 0.3                   | 14.16        |
| 22690 | 1.95      | 0.00         | 0.3                   | 9.75         |
| 22691 | 0.00      | 0.00         | 0.3                   | 9.30         |
| 22692 | 4.05      | 0.00         | 0.3                   | 24.35        |
| 22693 | 1.76      | 0.00         | 0.3                   | 10.56        |
| 22694 | 0.00      | 0.00         | 0.3                   | 5.80         |
| 22695 | 14.64     | 5.76         | 0.3                   | 73.20        |
| 22696 | 0.00      | 0.00         | 0.3                   | 5.30         |
| 22697 | 1.70      | 0.00         | 0.3                   | 13.00        |
| 22698 | 2.35      | 0.00         | 0.3                   | 14.15        |

[16]:
```python
#==> ENTER YOUR CODE HERE

df['payment_type'].unique()

# How many of each payment type are represented in the data?
```

```
[16]: array([1, 2, 3, 4])
```

According to the data dictionary, the payment method was encoded as follows:

1 = Credit card
2 = Cash
3 = No charge
4 = Dispute
5 = Unknown
6 = Voided trip

```
[23]: #==> ENTER YOUR CODE HERE

      # What is the average tip for trips paid for with credit card?
      df_cc_trips = df[df['payment_type'] == 1]
      print(df_cc_trips)

      #==> ENTER YOUR CODE HERE

      # What is the average tip for trips paid for with cash?
      print(df_cc_trips['tip_amount'].mean())
```

```
          Unnamed: 0  VendorID     tpep_pickup_datetime    tpep_dropoff_datetime  \
0           24870114         2   03/25/2017 8:55:43 AM    03/25/2017 9:09:47 AM
1           35634249         1   04/11/2017 2:53:28 PM    04/11/2017 3:19:58 PM
2          106203690         1   12/15/2017 7:26:56 AM    12/15/2017 7:34:08 AM
3           38942136         2   05/07/2017 1:17:59 PM    05/07/2017 1:48:14 PM
5           23345809         2   03/25/2017 8:34:11 PM    03/25/2017 8:42:11 PM
...              ...       ...                     ...                      ...
22692       60425673         1   07/16/2017 3:22:51 AM    07/16/2017 3:40:52 AM
22693       67858616         2  08/10/2017 10:20:04 PM   08/10/2017 10:29:31 PM
22695       66632549         2   08/06/2017 4:43:59 PM    08/06/2017 5:24:47 PM
22697       60217333         2  07/15/2017 12:56:30 PM    07/15/2017 1:08:26 PM
22698       17208911         1   03/02/2017 1:02:49 PM    03/02/2017 1:16:09 PM

          passenger_count  trip_distance  RatecodeID store_and_fwd_flag  \
0                       6           3.34           1                  N
1                       1           1.80           1                  N
2                       1           1.00           1                  N
3                       1           3.70           1                  N
5                       6           2.30           1                  N
...                   ...            ...         ...                ...
22692                   1           5.70           1                  N
22693                   1           0.89           1                  N
22695                   1          16.71           2                  N
22697                   1           2.36           1                  N
22698                   1           2.10           1                  N
```

```
       PULocationID   DOLocationID   payment_type   fare_amount   extra   mta_tax  \
0               100            231              1          13.0     0.0       0.5
1               186             43              1          16.0     0.0       0.5
2               262            236              1           6.5     0.0       0.5
3               188             97              1          20.5     0.0       0.5
5               161            236              1           9.0     0.5       0.5
...             ...            ...            ...           ...     ...       ...
22692           249             17              1          19.0     0.5       0.5
22693           229            170              1           7.5     0.5       0.5
22695           132            164              1          52.0     0.0       0.5
22697            68            144              1          10.5     0.0       0.5
22698           239            236              1          11.0     0.0       0.5

       tip_amount   tolls_amount   improvement_surcharge   total_amount
0            2.76           0.00                     0.3          16.56
1            4.00           0.00                     0.3          20.80
2            1.45           0.00                     0.3           8.75
3            6.39           0.00                     0.3          27.69
5            2.06           0.00                     0.3          12.36
...           ...            ...                     ...            ...
22692        4.05           0.00                     0.3          24.35
22693        1.76           0.00                     0.3          10.56
22695       14.64           5.76                     0.3          73.20
22697        1.70           0.00                     0.3          13.00
22698        2.35           0.00                     0.3          14.15

[15265 rows x 18 columns]
2.7298001965279934
```

```python
[29]: #==> ENTER YOUR CODE HERE

      # How many times is each vendor ID represented in the data?
      val_vendorIDs_types = df['VendorID'].value_counts()
      print(val_vendorIDs_types)
```

```
2    12626
1    10073
Name: VendorID, dtype: int64
```

```python
[37]: #==> ENTER YOUR CODE HERE

      # What is the mean total amount for each vendor?
      df_vendorID_1 = df[df['VendorID'] == 1]
      df_vendorID_2 = df[df['VendorID'] == 2]


      df_vendorID1_mean = df_vendorID_1['total_amount'].mean()
      df_vendorID2_mean = df_vendorID_2['total_amount'].mean()
```

```
print(df_vendorID1_mean)
print(df_vendorID2_mean)
```

```
16.298118733246966
16.32038175193886
```

[41]:
```
#==> ENTER YOUR CODE HERE

# Filter the data for credit card payments only
#done above and the dataframe is named as df_cc_trips

#==> ENTER YOUR CODE HERE

# Filter the credit-card-only data for passenger count only
df_passenger_count_cc = df_cc_trips['passenger_count']
df_passenger_count_cc.head(10)
```

[41]:
```
0     6
1     1
2     1
3     1
5     6
6     1
7     1
9     1
10    1
11    2
Name: passenger_count, dtype: int64
```

[56]:
```
#==> ENTER YOUR CODE HERE

# Calculate the average tip amount for each passenger count (credit card␣
 ↪payments only)

df_pcount_vals = df_cc_trips['passenger_count'].unique()
print(df_pcount_vals)

df_cc_trips.head()

df_cctrips_pcount1 = df_cc_trips[df_cc_trips['passenger_count']==1]
df_cctrips_pcount2 = df_cc_trips[df_cc_trips['passenger_count']==2]
df_cctrips_pcount3 = df_cc_trips[df_cc_trips['passenger_count']==3]
df_cctrips_pcount4 = df_cc_trips[df_cc_trips['passenger_count']==4]
df_cctrips_pcount5 = df_cc_trips[df_cc_trips['passenger_count']==5]
df_cctrips_pcount6 = df_cc_trips[df_cc_trips['passenger_count']==6]
```

```
#df_cctrips_pcount1.head()
#df_cctrips_pcount2.head()
#df_cctrips_pcount3.head()
#df_cctrips_pcount4.head()
#df_cctrips_pcount5.head()
#df_cctrips_pcount6.head()
```

[6 1 2 4 3 5 0]

[56]:

|    | Unnamed: 0 | VendorID | tpep_pickup_datetime | tpep_dropoff_datetime | \ |
|----|-----------|----------|----------------------|-----------------------|---|
| 11 | 18600059  | 2        | 03/05/2017 7:15:30 PM | 03/05/2017 7:52:18 PM |   |
| 17 | 58691513  | 2        | 07/10/2017 1:36:31 PM | 07/10/2017 1:48:43 PM |   |
| 19 | 18383214  | 2        | 03/05/2017 4:01:07 AM | 03/05/2017 4:14:11 AM |   |
| 43 | 34535675  | 2        | 04/26/2017 11:18:27 AM | 04/26/2017 11:24:59 AM |   |
| 57 | 19282330  | 2        | 03/07/2017 8:19:42 PM | 03/07/2017 8:29:08 PM |   |

|    | passenger_count | trip_distance | RatecodeID | store_and_fwd_flag | \ |
|----|-----------------|---------------|------------|--------------------|---|
| 11 | 2               | 18.90         | 2          | N                  |   |
| 17 | 2               | 1.71          | 1          | N                  |   |
| 19 | 2               | 2.77          | 1          | N                  |   |
| 43 | 2               | 1.12          | 1          | N                  |   |
| 57 | 2               | 2.51          | 1          | N                  |   |

|    | PULocationID | DOLocationID | payment_type | fare_amount | extra | mta_tax | \ |
|----|--------------|--------------|--------------|-------------|-------|---------|---|
| 11 | 236          | 132          | 1            | 52.0        | 0.0   | 0.5     |   |
| 17 | 142          | 100          | 1            | 9.5         | 0.0   | 0.5     |   |
| 19 | 79           | 68           | 1            | 11.5        | 0.5   | 0.5     |   |
| 43 | 234          | 211          | 1            | 6.5         | 0.0   | 0.5     |   |
| 57 | 229          | 146          | 1            | 10.0        | 0.5   | 0.5     |   |

|    | tip_amount | tolls_amount | improvement_surcharge | total_amount |
|----|------------|--------------|-----------------------|--------------|
| 11 | 14.58      | 5.54         | 0.3                   | 72.92        |
| 17 | 0.00       | 0.00         | 0.3                   | 10.30        |
| 19 | 3.20       | 0.00         | 0.3                   | 16.00        |
| 43 | 1.82       | 0.00         | 0.3                   | 9.12         |
| 57 | 3.39       | 0.00         | 0.3                   | 14.69        |

## 4.3 PACE: Construct

**Note**: The Construct stage does not apply to this workflow. The PACE framework can be adapted to fit the specific requirements of any project.

## 4.4 PACE: Execute

Consider the questions in your PACE Strategy Document and those below to craft your response.

### 4.4.1 Given your efforts, what can you summarize for DeShawn and the data team?

*Note for Learners: Your notebook should contain data that can address Luana's requests. Which two variables are most helpful for building a predictive model for the client: NYC TLC?*

==> ENTER YOUR RESPONSE HERE

**Congratulations!** You've completed this lab. However, you may not notice a green check mark next to this item on Coursera's platform. Please continue your progress regardless of the check mark. Just click on the "save" icon at the top of this notebook to ensure your work has been logged.