# Activity_Course 3 Automatidata project lab

February 4, 2024

## 1 Course 3 Automatidata project

**Course 3 - Go Beyond the Numbers: Translate Data into Insights**

You are the newest data professional in a fictional data consulting firm: Automatidata. The team is still early into the project, having only just completed an initial plan of action and some early Python coding work.

Luana Rodriquez, the senior data analyst at Automatidata, is pleased with the work you have already completed and requests your assistance with some EDA and data visualization work for the New York City Taxi and Limousine Commission project (New York City TLC) to get a general understanding of what taxi ridership looks like. The management team is asking for a Python notebook showing data structuring and cleaning, as well as any matplotlib/seaborn visualizations plotted to help understand the data. At the very least, include a box plot of the ride durations and some time series plots, like a breakdown by quarter or month.

Additionally, the management team has recently asked all EDA to include Tableau visualizations. For this taxi data, create a Tableau dashboard showing a New York City map of taxi/limo trips by month. Make sure it is easy to understand to someone who isn't data savvy, and remember that the assistant director at the New York City TLC is a person with visual impairments.

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

## 2 Course 3 End-of-course project: Exploratory data analysis

In this activity, you will examine data provided and prepare it for analysis. You will also design a professional data visualization that tells a story, and will help data-driven decisions for business needs.

Please note that the Tableau visualization activity is optional, and will not affect your completion of the course. Completing the Tableau activity will help you practice planning out and plotting a data visualization based on a specific business need. The structure of this activity is designed to emulate the proposals you will likely be assigned in your career as a data professional. Completing this activity will help prepare you for those career moments.

**The purpose** of this project is to conduct exploratory data analysis on a provided data set. Your mission is to continue the investigation you began in C2 and perform further EDA on this data with the aim of learning more about the variables.

**The goal** is to clean data set and create a visualization.
*This activity has 4 parts:*

**Part 1:** Imports, links, and loading

**Part 2:** Data Exploration * Data cleaning

**Part 3:** Building visualizations

**Part 4:** Evaluate and share results

Follow the instructions and answer the questions below 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 Visualize a story in Tableau and Python

# 4 PACE stages

- [Plan](#scrollTo=psz51YkZVwtN&line=3&uniqifier=1)
- [Analyze](#scrollTo=mA7Mz_SnI8km&line=4&uniqifier=1)
- [Construct](#scrollTo=Lca9c8XON8lc&line=2&uniqifier=1)
- [Execute](#scrollTo=401PgchTPr4E&line=2&uniqifier=1)

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

In this stage, consider the following questions where applicable to complete your code response: 1. Identify any outliers:

- What methods are best for identifying outliers?
- How do you make the decision to keep or exclude outliers from any future models?

==> ENTER YOUR RESPONSE HERE

### 4.1.1 Task 1. Imports, links, and loading

Go to Tableau Public The following link will help you complete this activity. Keep Tableau Public open as you proceed to the next steps.

Link to supporting materials: Tableau Public: https://public.tableau.com/s/

For EDA of the data, import the data and packages that would be most helpful, such as pandas, numpy and matplotlib.

```
[ ]: # Import packages and libraries
     #==> ENTER YOUR CODE HERE
```

**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.

```
[ ]: # Load dataset into dataframe
     df = pd.read_csv('2017_Yellow_Taxi_Trip_Data.csv')
```

## 4.2 PACE: Analyze

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

### 4.2.1 Task 2a. Data exploration and cleaning

Decide which columns are applicable

The first step is to assess your data. Check the Data Source page on Tableau Public to get a sense of the size, shape and makeup of the data set. Then answer these questions to yourself:

Given our scenario, which data columns are most applicable? Which data columns can I eliminate, knowing they won't solve our problem scenario?

Consider functions that help you understand and structure the data.

- head()
- describe()
- info()
- groupby()
- sortby()

What do you do about missing data (if any)?

Are there data outliers? What are they and how might you handle them?

What do the distributions of your variables tell you about the question you're asking or the problem you're trying to solve?

==> ENTER YOUR RESPONSE HERE

Start by discovering, using head and size.

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

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

Use describe...

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

And info.

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

### 4.2.2 Task 2b. Assess whether dimensions and measures are correct

On the data source page in Tableau, double check the data types for the applicable columns you selected on the previous step. Pay close attention to the dimensions and measures to assure they are correct.

In Python, consider the data types of the columns. *Consider:* Do they make sense?

Review the link provided in the previous activity instructions to create the required Tableau visualization.

### 4.2.3 Task 2c. Select visualization type(s)

Select data visualization types that will help you understand and explain the data.

Now that you know which data columns you'll use, it is time to decide which data visualization makes the most sense for EDA of the TLC dataset. What type of data visualization(s) would be most helpful?

- Line graph
- Bar chart
- Box plot
- Histogram
- Heat map
- Scatter plot
- A geographic map

==> ENTER YOUR RESPONSE HERE

## 4.3 PACE: Construct

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

### 4.3.1 Task 3. Data visualization

You've assessed your data, and decided on which data variables are most applicable. It's time to plot your visualization(s)!

### 4.3.2 Boxplots

Perform a check for outliers on relevant columns such as trip distance and trip duration. Remember, some of the best ways to identify the presence of outliers in data are box plots and histograms.

**Note:** Remember to convert your date columns to datetime in order to derive total trip duration.

```
[ ]: # Convert data columns to datetime
     #==> ENTER YOUR CODE HERE
```

**trip distance**

```
[ ]: # Create box plot of trip_distance
     #==> ENTER YOUR CODE HERE
```

```
[ ]: # Create histogram of trip_distance
     #==> ENTER YOUR CODE HERE
```

**total amount**

```
[ ]: # Create box plot of total_amount
     #==> ENTER YOUR CODE HERE
```

```
[ ]: # Create histogram of total_amount
     #==> ENTER YOUR CODE HERE
```

**tip amount**

```
[ ]: # Create box plot of tip_amount
     #==> ENTER YOUR CODE HERE
```

```
[ ]: # Create histogram of tip_amount
     #==> ENTER YOUR CODE HERE
```

**tip_amount by vendor**

```
[ ]: # Create histogram of tip_amount by vendor
     #==> ENTER YOUR CODE HERE
```

Next, zoom in on the upper end of the range of tips to check whether vendor one gets noticeably more of the most generous tips.

```
[ ]: # Create histogram of tip_amount by vendor for tips > $10
     #==> ENTER YOUR CODE HERE
```

**Mean tips by passenger count**

Examine the unique values in the `passenger_count` column.

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

```
[ ]: # Calculate mean tips by passenger_count
     #==> ENTER YOUR CODE HERE
```

```
[ ]: # Create bar plot for mean tips by passenger count
     #==> ENTER YOUR CODE HERE
```

**Create month and day columns**

```
[ ]: # Create a month column
     #==> ENTER YOUR CODE HERE

     # Create a day column
     #==> ENTER YOUR CODE HERE
```

**Plot total ride count by month**

Begin by calculating total ride count by month.

```
[ ]: # Get total number of rides for each month
     #==> ENTER YOUR CODE HERE
```

Reorder the results to put the months in calendar order.

```
[ ]: # Reorder the monthly ride list so months go in order
     #==> ENTER YOUR CODE HERE
```

```
[ ]: # Show the index
     #==> ENTER YOUR CODE HERE
```

```
[ ]: # Create a bar plot of total rides per month
     #==> ENTER YOUR CODE HERE
```

**Plot total ride count by day**

Repeat the above process, but now calculate the total rides by day of the week.

```
[ ]: # Repeat the above process, this time for rides by day
     #==> ENTER YOUR CODE HERE
```

```
[ ]: # Create bar plot for ride count by day
     #==> ENTER YOUR CODE HERE
```

**Plot total revenue by day of the week**

Repeat the above process, but now calculate the total revenue by day of the week.

```
[ ]: # Repeat the process, this time for total revenue by day
     #==> ENTER YOUR CODE HERE
```

```
[ ]: # Create bar plot of total revenue by day
     #==> ENTER YOUR CODE HERE
```

**Plot total revenue by month**

```
[ ]: # Repeat the process, this time for total revenue by month
     #==> ENTER YOUR CODE HERE
```

```
[ ]: # Create a bar plot of total revenue by month
     #==> ENTER YOUR CODE HERE
```

**Scatter plot**  You can create a scatterplot in Tableau Public, which can be easier to manipulate
and present. If you'd like step by step instructions, you can review the following link. Those
instructions create a scatterplot showing the relationship between total_amount and trip_distance.
Consider adding the Tableau visualization to your executive summary, and adding key insights from
your findings on those two variables.

Tableau visualization guidelines

**Plot mean trip distance by drop-off location**

```
[ ]: # Get number of unique drop-off location IDs
     #==> ENTER YOUR CODE HERE
```

```
[ ]: # Calculate the mean trip distance for each drop-off location
     #==> ENTER YOUR CODE HERE

     # Sort the results in descending order by mean trip distance
     #==> ENTER YOUR CODE HERE
```

```
[ ]: # Create a bar plot of mean trip distances by drop-off location in ascending␣
     ↪order by distance
     #==> ENTER YOUR CODE HERE
```

## 4.4  BONUS CONTENT

To confirm your conclusion, consider the following experiment: 1. Create a sample of coordinates
from a normal distribution—in this case 1,500 pairs of points from a normal distribution with a
mean of 10 and a standard deviation of 5 2. Calculate the distance between each pair of coordinates
3. Group the coordinates by endpoint and calculate the mean distance between that endpoint and
all other points it was paired with 4. Plot the mean distance for each unique endpoint

```
[ ]: #BONUS CONTENT

     #1. Generate random points on a 2D plane from a normal distribution
     #==> ENTER YOUR CODE HERE
```

```
# 2. Calculate Euclidean distances between points in first half and second half␣
 ↪of array
#==> ENTER YOUR CODE HERE

# 3. Group the coordinates by "drop-off location", compute mean distance
#==> ENTER YOUR CODE HERE

# 4. Plot the mean distance between each endpoint ("drop-off location") and all␣
 ↪points it connected to
#==> ENTER YOUR CODE HERE
```

**Histogram of rides by drop-off location**

First, check to whether the drop-off locations IDs are consecutively numbered. For instance, does it go 1, 2, 3, 4…, or are some numbers missing (e.g., 1, 3, 4…). If numbers aren't all consecutive, the histogram will look like some locations have very few or no rides when in reality there's no bar because there's no location.

```
[ ]: # Check if all drop-off locations are consecutively numbered
     #==> ENTER YOUR CODE HERE
```

To eliminate the spaces in the historgram that these missing numbers would create, sort the unique drop-off location values, then convert them to strings. This will make the histplot function display all bars directly next to each other.

```
[ ]: #==> ENTER YOUR CODE HERE
     # DOLocationID column is numeric, so sort in ascending order
     #==> ENTER YOUR CODE HERE

     # Convert to string
     #==> ENTER YOUR CODE HERE

     # Plot
     #==> ENTER YOUR CODE HERE
```

## 4.5   PACE: Execute

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

### 4.5.1   Task 4a. Results and evaluation

Having built visualizations in Tableau and in Python, what have you learned about the dataset? What other questions have your visualizations uncovered that you should pursue?

***Pro tip:*** Put yourself in your client's perspective, what would they want to know?

Use the following code fields to pursue any additional EDA based on the visualizations you've already plotted. Also use the space to make sure your visualizations are clean, easily understandable, and accessible.

***Ask yourself:*** Did you consider color, contrast, emphasis, and labeling?

==> ENTER YOUR RESPONSE HERE

I have learned ....

My other questions are ....

My client would likely want to know ...

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

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

### 4.5.2 Task 4b. Conclusion

*Make it professional and presentable*

You have visualized the data you need to share with the director now. Remember, the goal of a data visualization is for an audience member to glean the information on the chart in mere seconds.

*Questions to ask yourself for reflection:* Why is it important to conduct Exploratory Data Analysis? Why are the data visualizations provided in this notebook useful?

EDA is important because ... ==> ENTER YOUR RESPONSE HERE

Visualizations helped me understand .. ==> ENTER YOUR RESPONSE HERE

You've now completed professional data visualizations according to a business need. Well done!

**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.