# Activity_Course 3 Waze project lab

February 11, 2024

## 1 Waze Project

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

Your team is still in the early stages of their user churn project. So far, you've completed a project proposal and used Python to inspect and organize Waze's user data.

You check your inbox and notice a new message from Chidi Ga, your team's Senior Data Analyst. Chidi is pleased with the work you have already completed and requests your assistance with exploratory data analysis (EDA) and further data visualization. Harriet Hadzic, Waze's Director of Data Analysis, will want to review a Python notebook that shows your data exploration and visualization.

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

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

In this activity, you will examine data provided and prepare it for analysis.

**The purpose** of this project is to conduct exploratory data analysis (EDA) on a provided dataset.

**The goal** is to continue the examination of the data that you began in the previous Course, adding relevant visualizations that help communicate the story that the data tells.

*This activity has 4 parts:*

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

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

**Part 3:** Building visualizations

**Part 4:** Evaluating and sharing results

Follow the instructions and answer the question 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 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 to reflect on the Plan stage.

### 4.1.1 Task 1. Imports and data loading

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

```
[ ]: ### YOUR CODE HERE ###
```

Read in the data and store it as a dataframe object called df.

**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 the dataset into a dataframe
     df = pd.read_csv('waze_dataset.csv')
```

## 4.2 PACE: Analyze

Consider the questions in your PACE Strategy Document and those below where applicable to complete your code: 1. Does the data need to be restructured or converted into usable formats?

2. Are there any variables that have missing data?

==> ENTER YOUR RESPONSES TO QUESTIONS 1-2 HERE

### 4.2.1 Task 2. Data exploration and cleaning

Consider the following questions:

1. Given the scenario, which data columns are most applicable?

2. Which data columns can you eliminate, knowing they won't solve your problem scenario?

3. How would you check for missing data? And how would you handle missing data (if any)?

4. How would you check for outliers? And how would handle outliers (if any)?

==> ENTER YOUR RESPONSES TO QUESTIONS 1-4 HERE

**Data overview and summary statistics** Use the following methods and attributes on the dataframe:

- `head()`
- `size`
- `describe()`
- `info()`

It's always helpful to have this information at the beginning of a project, where you can always refer back to if needed.

```
[ ]:  ### YOUR CODE HERE ###
```

```
[ ]:  ### YOUR CODE HERE ###
```

Generate summary statistics using the `describe()` method.

```
[ ]:  ### YOUR CODE HERE ###
```

And summary information using the `info()` method.

```
[ ]:  ### YOUR CODE HERE ###
```

## 4.3 PACE: Construct

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

Consider the following questions as you prepare to deal with outliers:

1. What are some ways to identify outliers?
2. How do you make the decision to keep or exclude outliers from any future models?

==> ENTER YOUR RESPONSES TO QUESTIONS 1-2 HERE

### 4.3.1 Task 3a. Visualizations

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 Waze dataset.

**Question:** What type of data visualization(s) will be most helpful?

- Line graph
- Bar chart
- Box plot
- Histogram
- Heat map

- Scatter plot
- A geographic map

==> ENTER YOUR RESPONSE HERE

Begin by examining the spread and distribution of important variables using box plots and histograms.

**sessions** *The number of occurrence of a user opening the app during the month*

```
[ ]: # Box plot
     ### YOUR CODE HERE ###
```

```
[ ]: # Histogram
     ### YOUR CODE HERE ###
```

The `sessions` variable is a right-skewed distribution with half of the observations having 56 or fewer sessions. However, as indicated by the boxplot, some users have more than 700.

**drives** *An occurrence of driving at least 1 km during the month*

```
[ ]: # Box plot
     ### YOUR CODE HERE ###
```

```
[ ]: # Histogram
     ### YOUR CODE HERE ###
```

The `drives` information follows a distribution similar to the `sessions` variable. It is right-skewed, approximately log-normal, with a median of 48. However, some drivers had over 400 drives in the last month.

**total_sessions** *A model estimate of the total number of sessions since a user has onboarded*

```
[ ]: # Box plot
     ### YOUR CODE HERE ###
```

```
[ ]: # Histogram
     ### YOUR CODE HERE ###
```

The `total_sessions` is a right-skewed distribution. The median total number of sessions is 159.6. This is interesting information because, if the median number of sessions in the last month was 48 and the median total sessions was ~160, then it seems that a large proportion of a user's total drives might have taken place in the last month. This is something you can examine more closely later.

**n_days_after_onboarding** *The number of days since a user signed up for the app*

4

```
[ ]: # Box plot
     ### YOUR CODE HERE ###
```

```
[ ]: # Histogram
     ### YOUR CODE HERE ###
```

The total user tenure (i.e., number of days since onboarding) is a uniform distribution with values ranging from near-zero to ~3,500 (~9.5 years).

**driven_km_drives**    *Total kilometers driven during the month*

```
[ ]: # Box plot
     ### YOUR CODE HERE ###
```

```
[ ]: # Histogram
     ### YOUR CODE HERE ###
```

The number of drives driven in the last month per user is a right-skewed distribution with half the users driving under 3,495 kilometers. As you discovered in the analysis from the previous course, the users in this dataset drive *a lot*. The longest distance driven in the month was over half the circumferene of the earth.

**duration_minutes_drives**    *Total duration driven in minutes during the month*

```
[ ]: # Box plot
     ### YOUR CODE HERE ###
```

```
[ ]: # Histogram
     ### YOUR CODE HERE ###
```

The `duration_minutes_drives` variable has a heavily skewed right tail. Half of the users drove less than ~1,478 minutes (~25 hours), but some users clocked over 250 hours over the month.

**activity_days**    *Number of days the user opens the app during the month*

```
[ ]: # Box plot
     ### YOUR CODE HERE ###
```

```
[ ]: # Histogram
     ### YOUR CODE HERE ###
```

Within the last month, users opened the app a median of 16 times. The box plot reveals a centered distribution. The histogram shows a nearly uniform distribution of ~500 people opening the app on each count of days. However, there are ~250 people who didn't open the app at all and ~250 people who opened the app every day of the month.

This distribution is noteworthy because it does not mirror the `sessions` distribution, which you might think would be closely correlated with `activity_days`.

**driving_days**  *Number of days the user drives (at least 1 km) during the month*

```
[ ]: # Box plot
     ### YOUR CODE HERE ###
```

```
[ ]: # Histogram
     ### YOUR CODE HERE ###
```

The number of days users drove each month is almost uniform, and it largely correlates with the number of days they opened the app that month, except the `driving_days` distribution tails off on the right.

However, there were almost twice as many users (~1,000 vs. ~550) who did not drive at all during the month. This might seem counterintuitive when considered together with the information from `activity_days`. That variable had ~500 users opening the app on each of most of the day counts, but there were only ~250 users who did not open the app at all during the month and ~250 users who opened the app every day. Flag this for further investigation later.

**device**  *The type of device a user starts a session with*

This is a categorical variable, so you do not plot a box plot for it. A good plot for a binary categorical variable is a pie chart.

```
[ ]: # Pie chart
     ### YOUR CODE HERE ###
```

There are nearly twice as many iPhone users as Android users represented in this data.

**label**  *Binary target variable ("retained" vs "churned") for if a user has churned anytime during the course of the month*

This is also a categorical variable, and as such would not be plotted as a box plot. Plot a pie chart instead.

```
[ ]: # Pie chart
     ### YOUR CODE HERE ###
```

Less than 18% of the users churned.

**driving_days vs. activity_days**  Because both `driving_days` and `activity_days` represent counts of days over a month and they're also closely related, you can plot them together on a single histogram. This will help to better understand how they relate to each other without having to scroll back and forth comparing histograms in two different places.

Plot a histogram that, for each day, has a bar representing the counts of `driving_days` and `activity_days`.

```
[ ]:  # Histogram
      ### YOUR CODE HERE ###
```

As observed previously, this might seem counterintuitive. After all, why are there *fewer* people who didn't use the app at all during the month and *more* people who didn't drive at all during the month?

On the other hand, it could just be illustrative of the fact that, while these variables are related to each other, they're not the same. People probably just open the app more than they use the app to drive—perhaps to check drive times or route information, to update settings, or even just by mistake.

Nonetheless, it might be worthwile to contact the data team at Waze to get more information about this, especially because it seems that the number of days in the month is not the same between variables.

Confirm the maximum number of days for each variable—`driving_days` and `activity_days`.

```
[ ]:  ### YOUR CODE HERE ###
```

It's true. Although it's possible that not a single user drove all 31 days of the month, it's highly unlikely, considering there are 15,000 people represented in the dataset.

One other way to check the validity of these variables is to plot a simple scatter plot with the x-axis representing one variable and the y-axis representing the other.

```
[ ]:  # Scatter plot
      ### YOUR CODE HERE ###
```

Notice that there is a theoretical limit. If you use the app to drive, then by definition it must count as a day-use as well. In other words, you cannot have more drive-days than activity-days. None of the samples in this data violate this rule, which is good.

**Retention by device** Plot a histogram that has four bars—one for each device-label combination—to show how many iPhone users were retained/churned and how many Android users were retained/churned.

```
[ ]:  # Histogram
      ### YOUR CODE HERE ###
```

The proportion of churned users to retained users is consistent between device types.

**Retention by kilometers driven per driving day** In the previous course, you discovered that the median distance driven last month for users who churned was 8.33 km, versus 3.36 km for people who did not churn. Examine this further.

1. Create a new column in `df` called `km_per_driving_day`, which represents the mean distance driven per driving day for each user.

2. Call the `describe()` method on the new column.

```
[ ]: # 1. Create `km_per_driving_day` column
     ### YOUR CODE HERE ###

     # 2. Call `describe()` on the new column
     ### YOUR CODE HERE ###
```

What do you notice? The mean value is infinity, the standard deviation is NaN, and the max value is infinity. Why do you think this is?

This is the result of there being values of zero in the `driving_days` column. Pandas imputes a value of infinity in the corresponding rows of the new column because division by zero is undefined.

1. Convert these values from infinity to zero. You can use `np.inf` to refer to a value of infinity.

2. Call `describe()` on the `km_per_driving_day` column to verify that it worked.

```
[ ]: # 1. Convert infinite values to zero
     ### YOUR CODE HERE ###

     # 2. Confirm that it worked
     ### YOUR CODE HERE ###
```

The maximum value is 15,420 kilometers *per drive day*. This is physically impossible. Driving 100 km/hour for 12 hours is 1,200 km. It's unlikely many people averaged more than this each day they drove, so, for now, disregard rows where the distance in this column is greater than 1,200 km.

Plot a histogram of the new `km_per_driving_day` column, disregarding those users with values greater than 1,200 km. Each bar should be the same length and have two colors, one color representing the percent of the users in that bar that churned and the other representing the percent that were retained. This can be done by setting the `multiple` parameter of seaborn's `histplot()` function to `fill`.

```
[ ]: # Histogram
     ### YOUR CODE HERE ###
```

The churn rate tends to increase as the mean daily distance driven increases, confirming what was found in the previous course. It would be worth investigating further the reasons for long-distance users to discontinue using the app.

**Churn rate per number of driving days**   Create another histogram just like the previous one, only this time it should represent the churn rate for each number of driving days.

```
[ ]: # Histogram
     ### YOUR CODE HERE ###
```

The churn rate is highest for people who didn't use Waze much during the last month. The more times they used the app, the less likely they were to churn. While 40% of the users who didn't use the app at all last month churned, nobody who used the app 30 days churned.

This isn't surprising. If people who used the app a lot churned, it would likely indicate dissatisfaction. When people who don't use the app churn, it might be the result of dissatisfaction in the past, or it might be indicative of a lesser need for a navigational app. Maybe they moved to a city with good public transportation and don't need to drive anymore.

**Proportion of sessions that occurred in the last month** Create a new column `percent_sessions_in_last_month` that represents the percentage of each user's total sessions that were logged in their last month of use.

```
[ ]: ### YOUR CODE HERE ###
```

What is the median value of the new column?

```
[ ]: ### YOUR CODE HERE ###
```

Now, create a histogram depicting the distribution of values in this new column.

```
[ ]: # Histogram
     ### YOUR CODE HERE ###
```

Check the median value of the `n_days_after_onboarding` variable.

```
[ ]: ### YOUR CODE HERE ###
```

Half of the people in the dataset had 40% or more of their sessions in just the last month, yet the overall median time since onboarding is almost five years.

Make a histogram of `n_days_after_onboarding` for just the people who had 40% or more of their total sessions in the last month.

```
[ ]: # Histogram
     ### YOUR CODE HERE ###
```

The number of days since onboarding for users with 40% or more of their total sessions occurring in just the last month is a uniform distribution. This is very strange. It's worth asking Waze why so many long-time users suddenly used the app so much in the last month.

### 4.3.2 Task 3b. Handling outliers

The box plots from the previous section indicated that many of these variables have outliers. These outliers do not seem to be data entry errors; they are present because of the right-skewed distributions.

Depending on what you'll be doing with this data, it may be useful to impute outlying data with more reasonable values. One way of performing this imputation is to set a threshold based on a percentile of the distribution.

To practice this technique, write a function that calculates the 95th percentile of a given column, then imputes values > the 95th percentile with the value at the 95th percentile. such as the 95th percentile of the distribution.

```
[ ]: ### YOUR CODE HERE ###
```

Next, apply that function to the following columns: * `sessions` * `drives` * `total_sessions` * `driven_km_drives` * `duration_minutes_drives`

```
[ ]: ### YOUR CODE HERE ###
```

Call `describe()` to see if your change worked.

```
[ ]: ### YOUR CODE HERE ###
```

**Conclusion**  Analysis revealed that the overall churn rate is ~17%, and that this rate is consistent between iPhone users and Android users.

Perhaps you feel that the more deeply you explore the data, the more questions arise. This is not uncommon! In this case, it's worth asking the Waze data team why so many users used the app so much in just the last month.

Also, EDA has revealed that users who drive very long distances on their driving days are *more* likely to churn, but users who drive more often are *less* likely to churn. The reason for this discrepancy is an opportunity for further investigation, and it would be something else to ask the Waze data team about.

## 4.4  PACE: Execute

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

### 4.4.1  Task 4a. Results and evaluation

Having built visualizations 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 …

Use the following two code blocks (add more blocks if you like) to do additional EDA you feel is important based on the given scenario.

```
[ ]: ### YOUR CODE HERE ###
```

```
[ ]: ### YOUR CODE HERE ###
```

### 4.4.2  Task 4b.  Conclusion

Now that you've explored and visualized your data, the next step is to share your findings with Harriet Hadzic, Waze's Director of Data Analysis. Consider the following questions as you prepare to write your executive summary. Think about key points you may want to share with the team, and what information is most relevant to the user churn project.

**Questions:**

1. What types of distributions did you notice in the variables? What did this tell you about the data?

2. Was there anything that led you to believe the data was erroneous or problematic in any way?

3. Did your investigation give rise to further questions that you would like to explore or ask the Waze team about?

4. What percentage of users churned and what percentage were retained?

5. What factors correlated with user churn? How?

6. Did newer uses have greater representation in this dataset than users with longer tenure? How do you know?

==> ENTER YOUR RESPONSES TO QUESTIONS 1-6 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.