

Activity_Course 3 Waze project lab

October 27, 2025

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.

```
[3]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
```

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.

```
[4]: # 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?

1. The data is mostly well-structured. Numeric columns are ready for analysis. The 'label' and 'device' columns are categorical and will need to be encoded into numeric format for further analysis.
2. The 'label' column has 70 missing values. All other variables are complete. Rows with missing labels should be dropped since they cannot contribute to churn prediction.

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)?

Q1: The columns most applicable to churn analysis are sessions, drives, total_sessions, and n_days_after_onboarding, as they describe user activity and engagement.

Q2: Columns like user_id or other identifiers can be eliminated because they don't help predict churn.

Q3: I would use `df.isnull().sum()` to identify missing values. Missing numeric data would be filled with the median; missing categorical data would be filled with the mode.

Q4: I would detect outliers using boxplots and summary statistics. Outliers would be handled by removing values outside $1.5 \times \text{IQR}$, depending on their influence on the data distribution.

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.

```
[5]: df.head()
```

```
[5]:   ID  label  sessions  drives  total_sessions  n_days_after_onboarding  \
0   0  retained      283     226      296.748273                2276
1   1  retained      133     107      326.896596                1225
2   2  retained      114      95      135.522926                2651
3   3  retained       49      40       67.589221                 15
4   4  retained       84      68      168.247020               1562

      total_navigations_fav1  total_navigations_fav2  driven_km_drives  \
0                        208                        0      2628.845068
1                         19                       64     13715.920550
2                         0                        0      3059.148818
3                        322                        7       913.591123
4                        166                        5     3950.202008

      duration_minutes_drives  activity_days  driving_days  device
```

0	1985.775061	28	19	Android
1	3160.472914	13	11	iPhone
2	1610.735904	14	8	Android
3	587.196542	7	3	iPhone
4	1219.555924	27	18	Android

```
[6]: df.size
```

```
[6]: 194987
```

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

```
[7]: df.describe()
```

```
[7]:
```

	ID	sessions	drives	total_sessions	\
count	14999.000000	14999.000000	14999.000000	14999.000000	
mean	7499.000000	80.633776	67.281152	189.964447	
std	4329.982679	80.699065	65.913872	136.405128	
min	0.000000	0.000000	0.000000	0.220211	
25%	3749.500000	23.000000	20.000000	90.661156	
50%	7499.000000	56.000000	48.000000	159.568115	
75%	11248.500000	112.000000	93.000000	254.192341	
max	14998.000000	743.000000	596.000000	1216.154633	

	n_days_after_onboarding	total_navigations_fav1	\
count	14999.000000	14999.000000	
mean	1749.837789	121.605974	
std	1008.513876	148.121544	
min	4.000000	0.000000	
25%	878.000000	9.000000	
50%	1741.000000	71.000000	
75%	2623.500000	178.000000	
max	3500.000000	1236.000000	

	total_navigations_fav2	driven_km_drives	duration_minutes_drives	\
count	14999.000000	14999.000000	14999.000000	
mean	29.672512	4039.340921	1860.976012	
std	45.394651	2502.149334	1446.702288	
min	0.000000	60.441250	18.282082	
25%	0.000000	2212.600607	835.996260	
50%	9.000000	3493.858085	1478.249859	
75%	43.000000	5289.861262	2464.362632	
max	415.000000	21183.401890	15851.727160	

	activity_days	driving_days
count	14999.000000	14999.000000
mean	15.537102	12.179879

std	9.004655	7.824036
min	0.000000	0.000000
25%	8.000000	5.000000
50%	16.000000	12.000000
75%	23.000000	19.000000
max	31.000000	30.000000

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

```
[8]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 14999 entries, 0 to 14998
Data columns (total 13 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   ID                                     14999 non-null  int64
1   label                                14299 non-null  object
2   sessions                             14999 non-null  int64
3   drives                               14999 non-null  int64
4   total_sessions                       14999 non-null  float64
5   n_days_after_onboarding              14999 non-null  int64
6   total_navigations_fav1               14999 non-null  int64
7   total_navigations_fav2               14999 non-null  int64
8   driven_km_drives                     14999 non-null  float64
9   duration_minutes_drives              14999 non-null  float64
10  activity_days                         14999 non-null  int64
11  driving_days                          14999 non-null  int64
12  device                                14999 non-null  object
dtypes: float64(3), int64(8), object(2)
memory usage: 1.5+ MB
```

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?
1. Outliers can be identified using summary statistics, box plots, scatter plots, or by calculating Z-scores. These methods help visualize or quantify values that fall far outside the normal data range.
 2. The decision to keep or remove outliers depends on their cause and impact. If outliers represent real user behavior, they should be kept. If they result from errors or data inconsistencies, they should be removed or adjusted. For modeling, removing extreme outliers may improve performance and accuracy.

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?

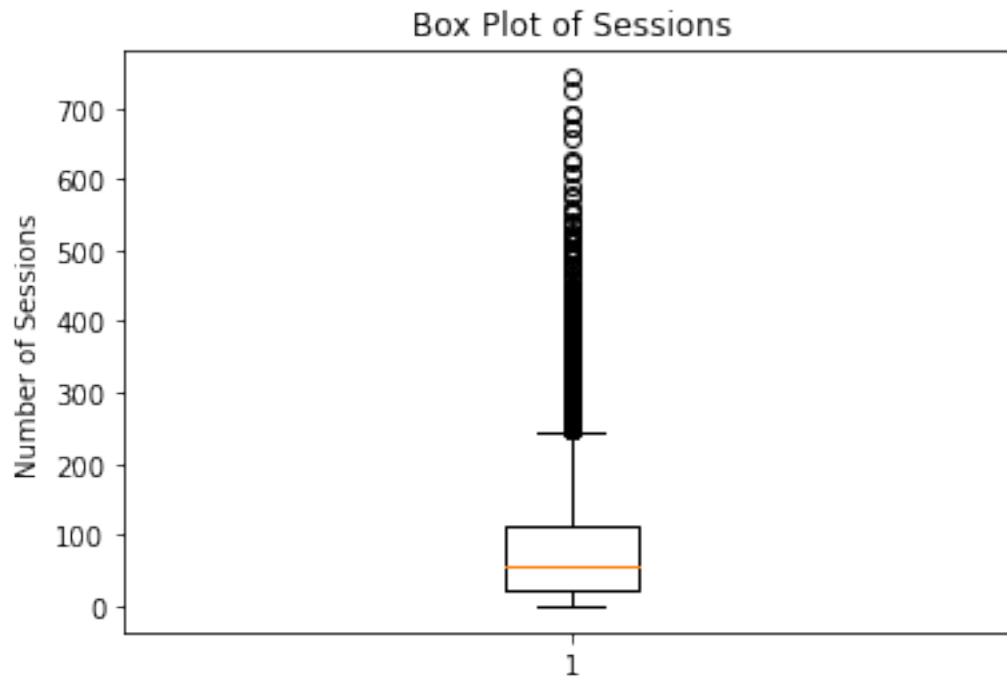
- Line graph
- Bar chart
- Box plot
- Histogram
- Heat map
- Scatter plot
- A geographic map

Box plots, histograms, scatter plots, and heat maps will be the most helpful visualizations for this dataset. Box plots and histograms will show the distribution and presence of outliers, while scatter plots will highlight relationships between key variables like drives and sessions. Bar charts can also help compare categorical variables such as device type and churn status.

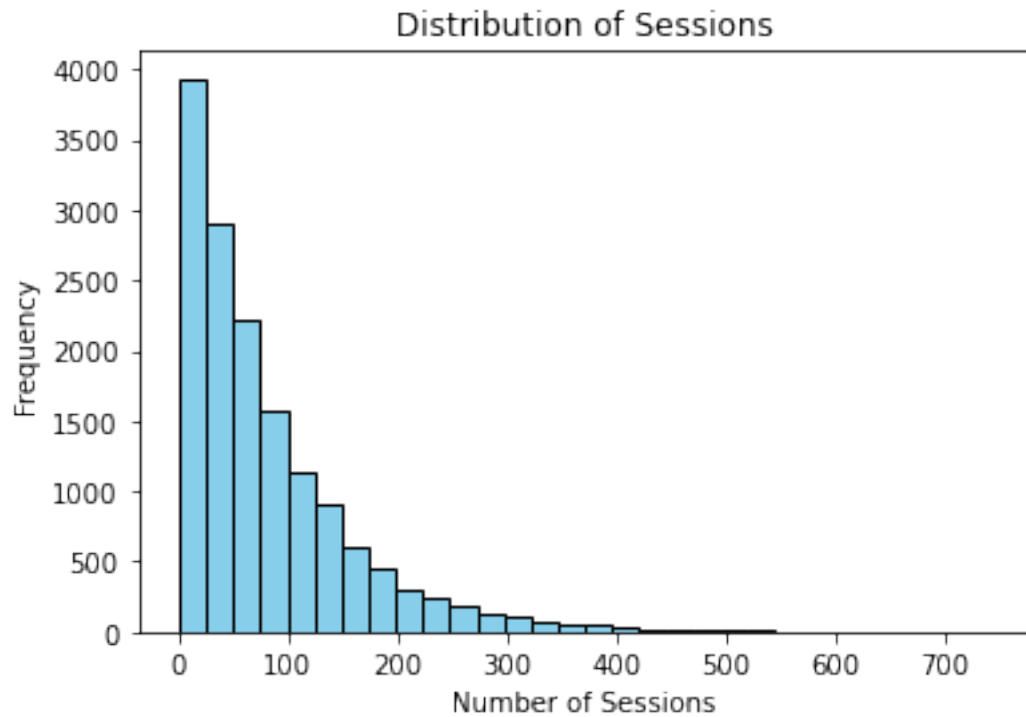
Begin by examining the spread and distribution of important variables using box plots and histograms.

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

```
[9]: # Box plot for 'sessions'
plt.figure(figsize=(6, 4))
plt.boxplot(df['sessions'])
plt.title('Box Plot of Sessions')
plt.ylabel('Number of Sessions')
plt.show()
```



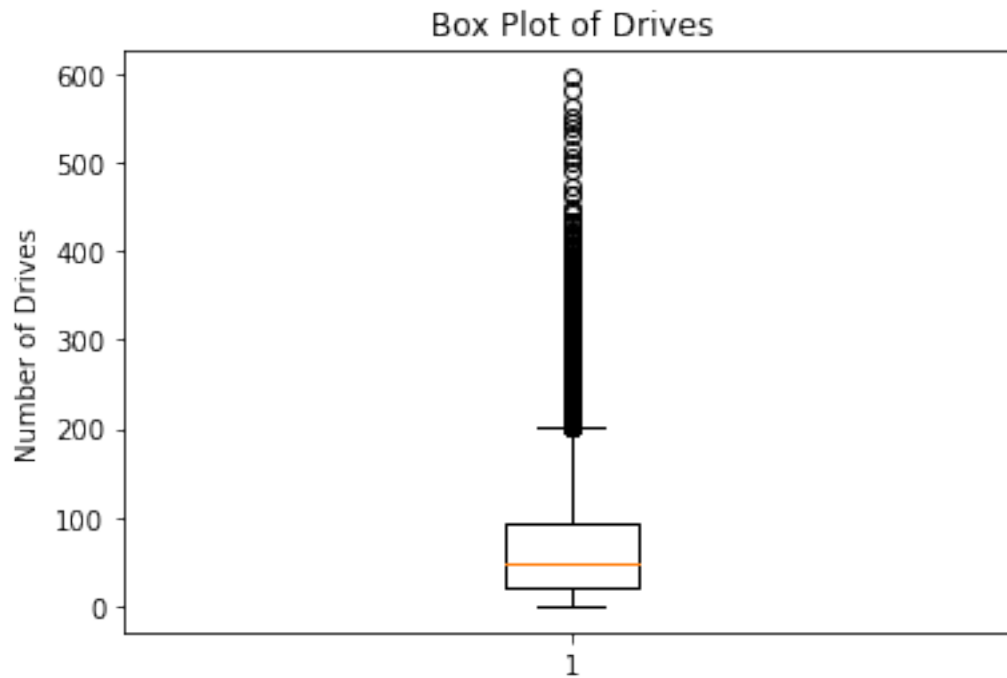
```
[10]: # Histogram for 'sessions'
plt.figure(figsize=(6, 4))
plt.hist(df['sessions'], bins=30, color='skyblue', edgecolor='black')
plt.title('Distribution of Sessions')
plt.xlabel('Number of Sessions')
plt.ylabel('Frequency')
plt.show()
```



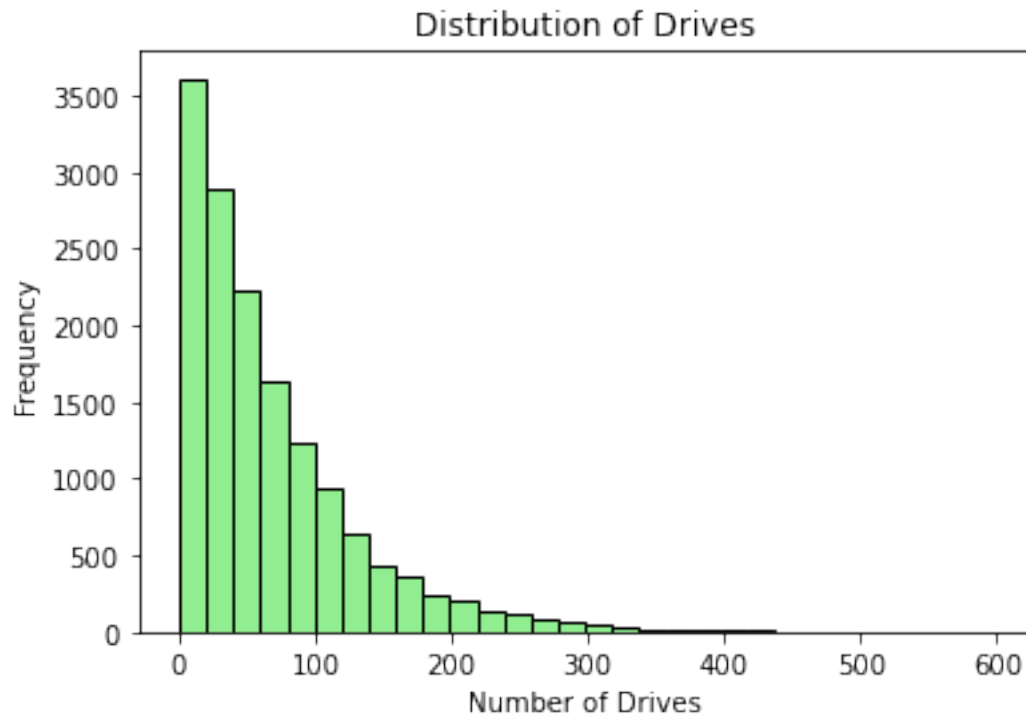
The `sessions` variable is a right-skewed distribution with half of the observations having 56 or fewer sessions. However, as indicated by the boxplot, some users have more than 700.

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

```
[11]: # Box plot for 'drives'
plt.figure(figsize=(6, 4))
plt.boxplot(df['drives'])
plt.title('Box Plot of Drives')
plt.ylabel('Number of Drives')
plt.show()
```

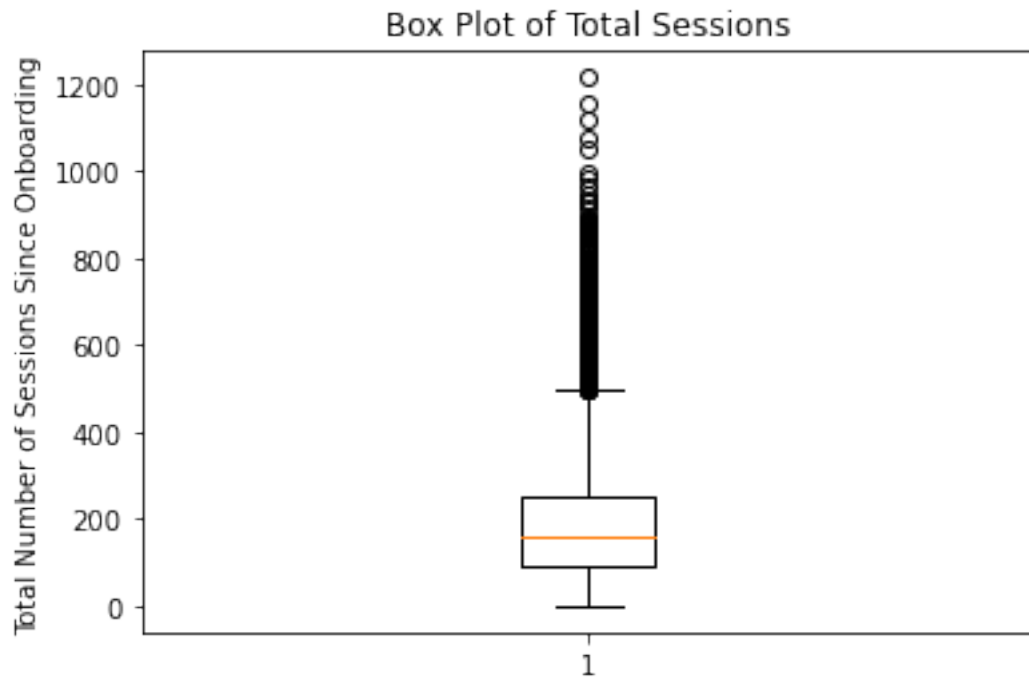
```
[12]: # Histogram for 'drives'
plt.figure(figsize=(6, 4))
plt.hist(df['drives'], bins=30, color='lightgreen', edgecolor='black')
plt.title('Distribution of Drives')
plt.xlabel('Number of Drives')
plt.ylabel('Frequency')
plt.show()
```



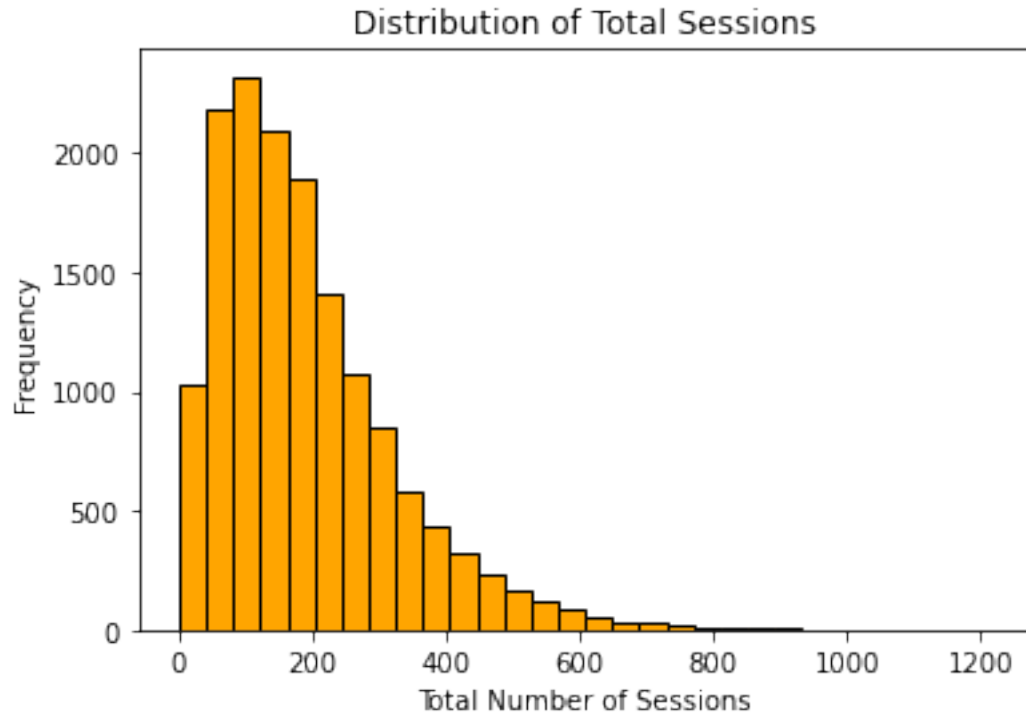
The **drives** information follows a distribution similar to the **sessions** variable. It is right-skewed, approximately log-normal, with a median of 48. However, some drivers had over 400 drives in the last month.

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

```
[13]: # Box plot for 'total_sessions'
plt.figure(figsize=(6, 4))
plt.boxplot(df['total_sessions'])
plt.title('Box Plot of Total Sessions')
plt.ylabel('Total Number of Sessions Since Onboarding')
plt.show()
```



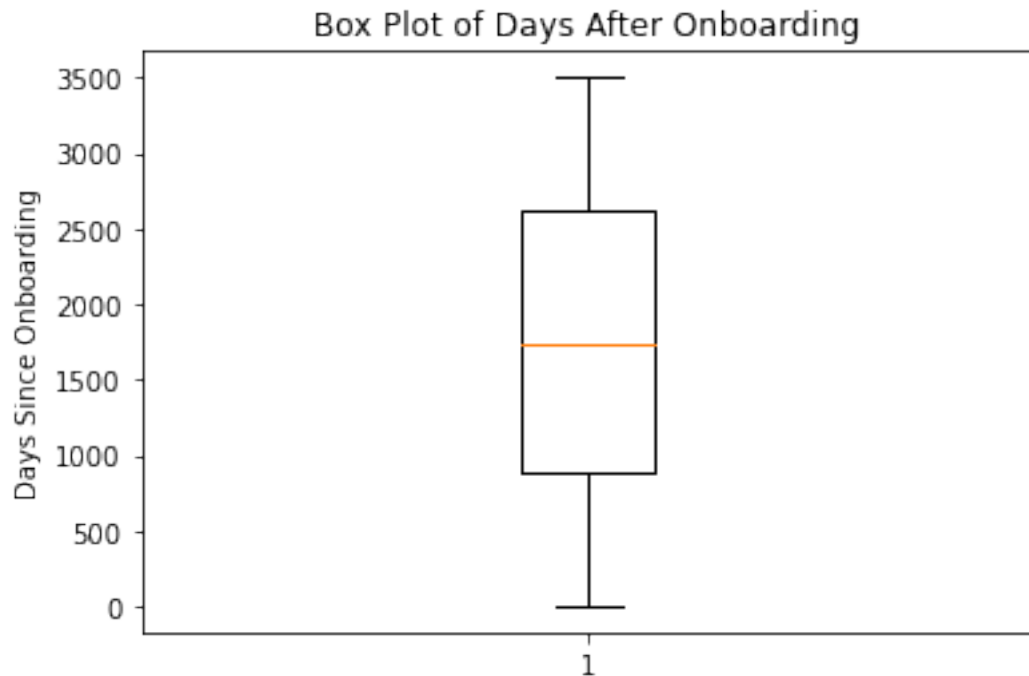
```
[14]: # Histogram for 'total_sessions'
plt.figure(figsize=(6, 4))
plt.hist(df['total_sessions'], bins=30, color='orange', edgecolor='black')
plt.title('Distribution of Total Sessions')
plt.xlabel('Total Number of Sessions')
plt.ylabel('Frequency')
plt.show()
```



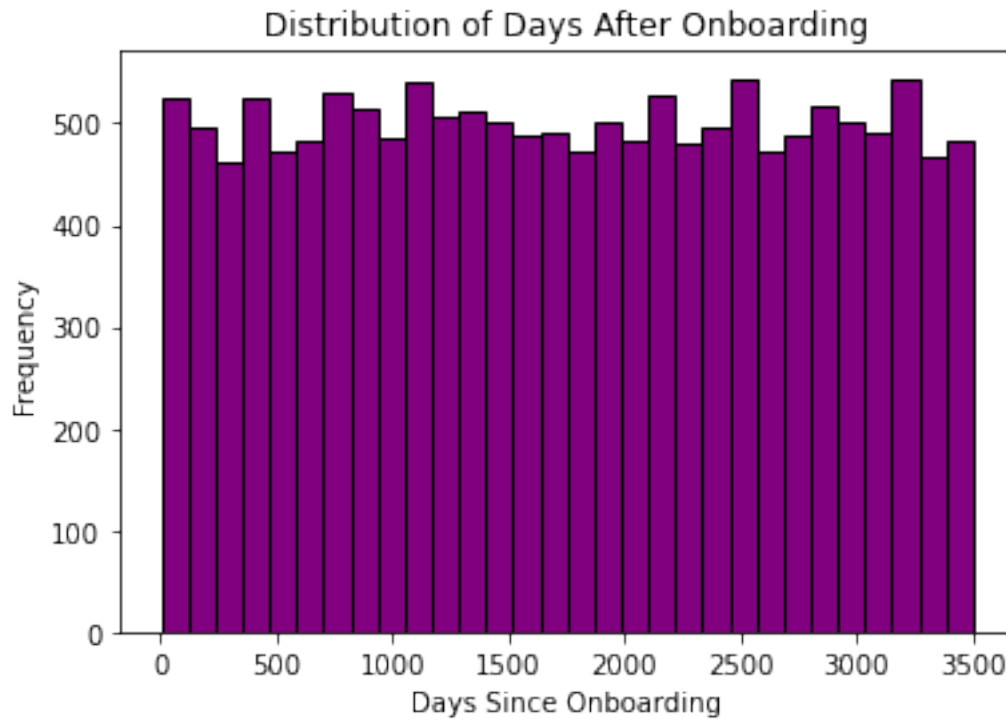
The `total_sessions` is a right-skewed distribution. The median total number of sessions is 159.6. This is interesting information because, if the median number of sessions in the last month was 48 and the median total sessions was ~160, then it seems that a large proportion of a user's total drives 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*

```
[15]: # Box plot for 'n_days_after_onboarding'
plt.figure(figsize=(6, 4))
plt.boxplot(df['n_days_after_onboarding'])
plt.title('Box Plot of Days After Onboarding')
plt.ylabel('Days Since Onboarding')
plt.show()
```



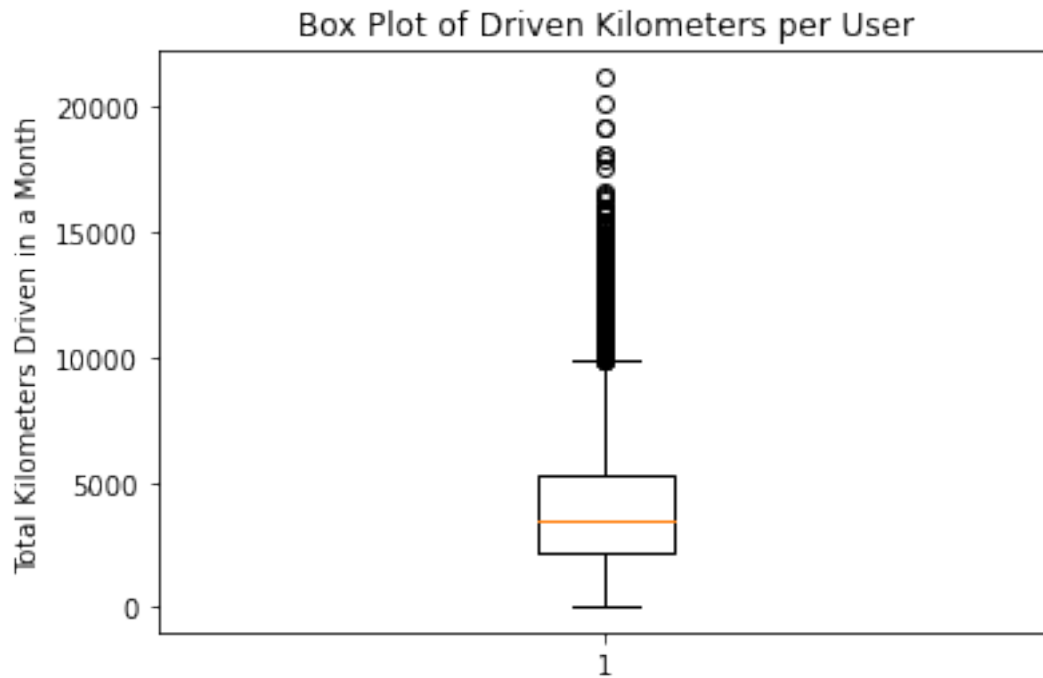
```
[16]: # Histogram for 'n_days_after_onboarding'
plt.figure(figsize=(6, 4))
plt.hist(df['n_days_after_onboarding'], bins=30, color='purple',
        edgecolor='black')
plt.title('Distribution of Days After Onboarding')
plt.xlabel('Days Since Onboarding')
plt.ylabel('Frequency')
plt.show()
```



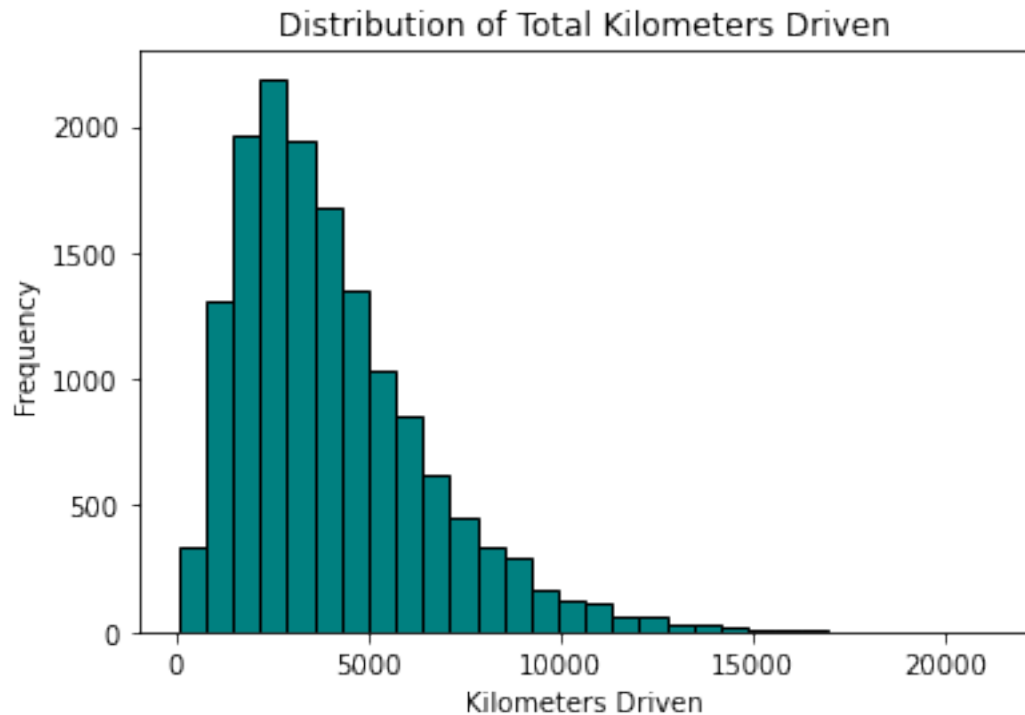
The total user tenure (i.e., number of days since onboarding) is a uniform distribution with values ranging from near-zero to ~3,500 (~9.5 years).

driven_km_drives *Total kilometers driven during the month*

```
[17]: # Box plot for 'driven_km_drives'
plt.figure(figsize=(6, 4))
plt.boxplot(df['driven_km_drives'])
plt.title('Box Plot of Driven Kilometers per User')
plt.ylabel('Total Kilometers Driven in a Month')
plt.show()
```



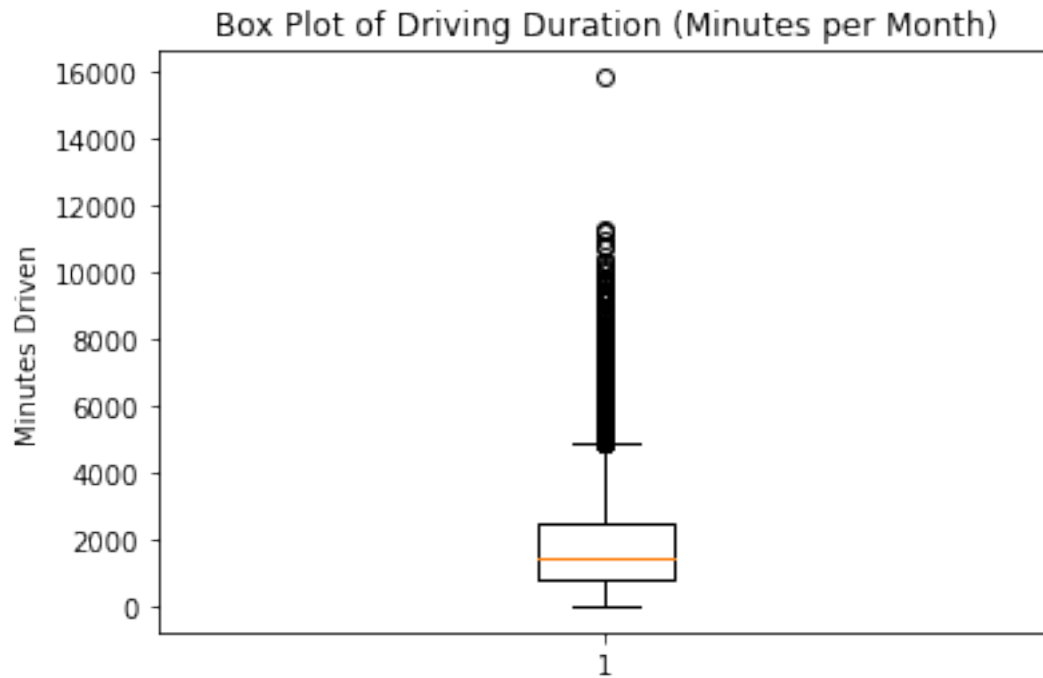
```
[18]: # Histogram for 'driven_km_drives'
plt.figure(figsize=(6, 4))
plt.hist(df['driven_km_drives'], bins=30, color='teal', edgecolor='black')
plt.title('Distribution of Total Kilometers Driven')
plt.xlabel('Kilometers Driven')
plt.ylabel('Frequency')
plt.show()
```



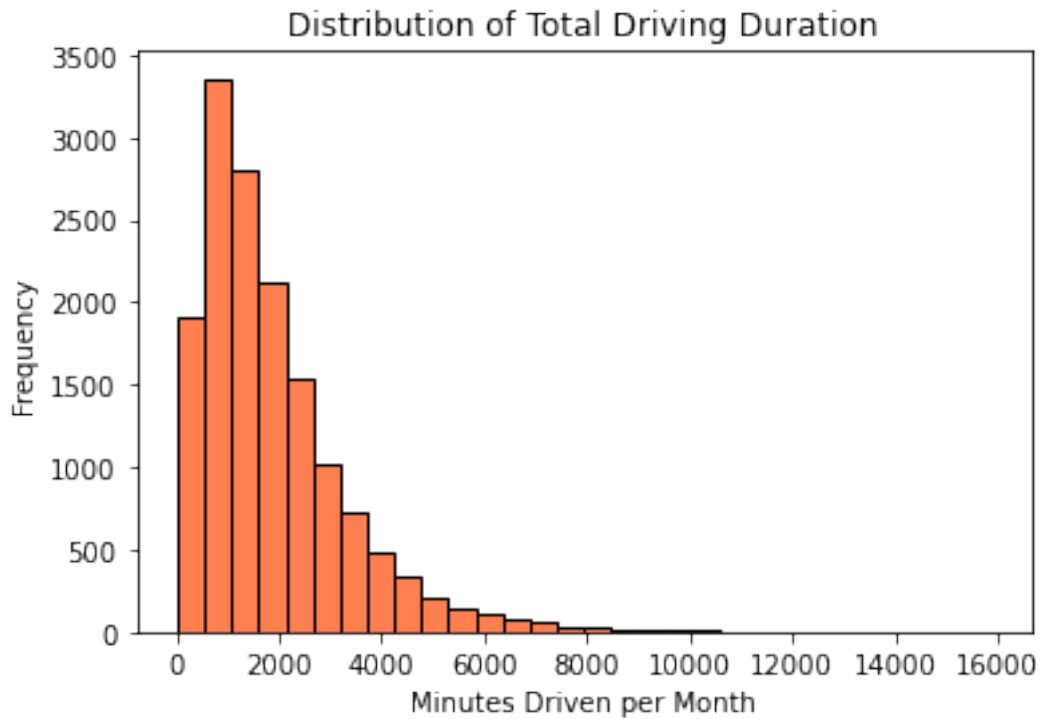
The number of drives driven in the last month per user is a right-skewed distribution with half the users driving under 3,495 kilometers. As you discovered in the analysis from the previous course, the users in this dataset drive *a lot*. The longest distance driven in the month was over half the circumference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

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

```
[19]: # Box plot for 'duration_minutes_drives'
plt.figure(figsize=(6, 4))
plt.boxplot(df['duration_minutes_drives'])
plt.title('Box Plot of Driving Duration (Minutes per Month)')
plt.ylabel('Minutes Driven')
plt.show()
```

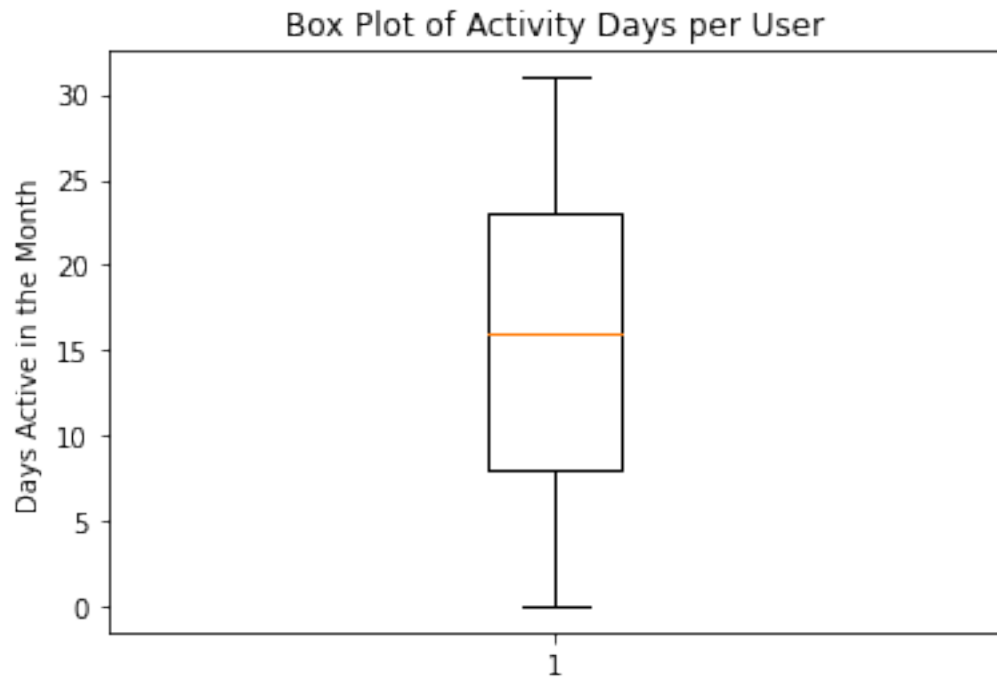
```
[20]: # Histogram for 'duration_minutes_drives'
plt.figure(figsize=(6, 4))
plt.hist(df['duration_minutes_drives'], bins=30, color='coral',
        edgecolor='black')
plt.title('Distribution of Total Driving Duration')
plt.xlabel('Minutes Driven per Month')
plt.ylabel('Frequency')
plt.show()
```



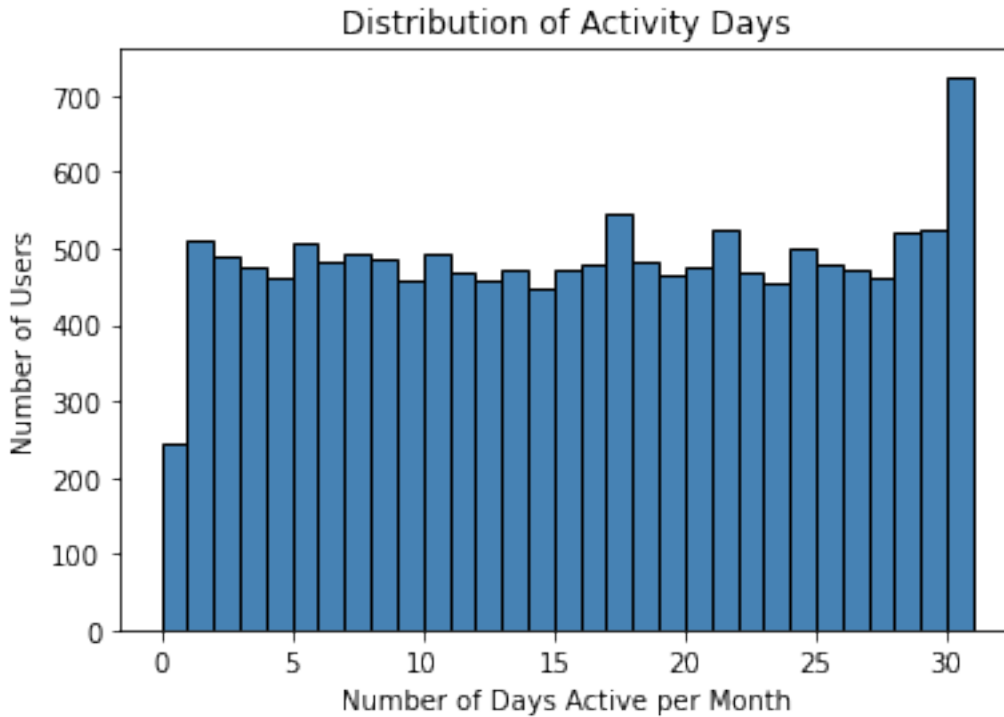
The `duration_minutes_drives` variable has a heavily skewed right tail. Half of the users drove less than ~1,478 minutes (~25 hours), but some users clocked over 250 hours over the month.

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

```
[21]: # Box plot for 'activity_days'
plt.figure(figsize=(6, 4))
plt.boxplot(df['activity_days'])
plt.title('Box Plot of Activity Days per User')
plt.ylabel('Days Active in the Month')
plt.show()
```



```
[22]: # Histogram for 'activity_days'
plt.figure(figsize=(6, 4))
plt.hist(df['activity_days'], bins=31, color='steelblue', edgecolor='black')
plt.title('Distribution of Activity Days')
plt.xlabel('Number of Days Active per Month')
plt.ylabel('Number of Users')
plt.show()
```

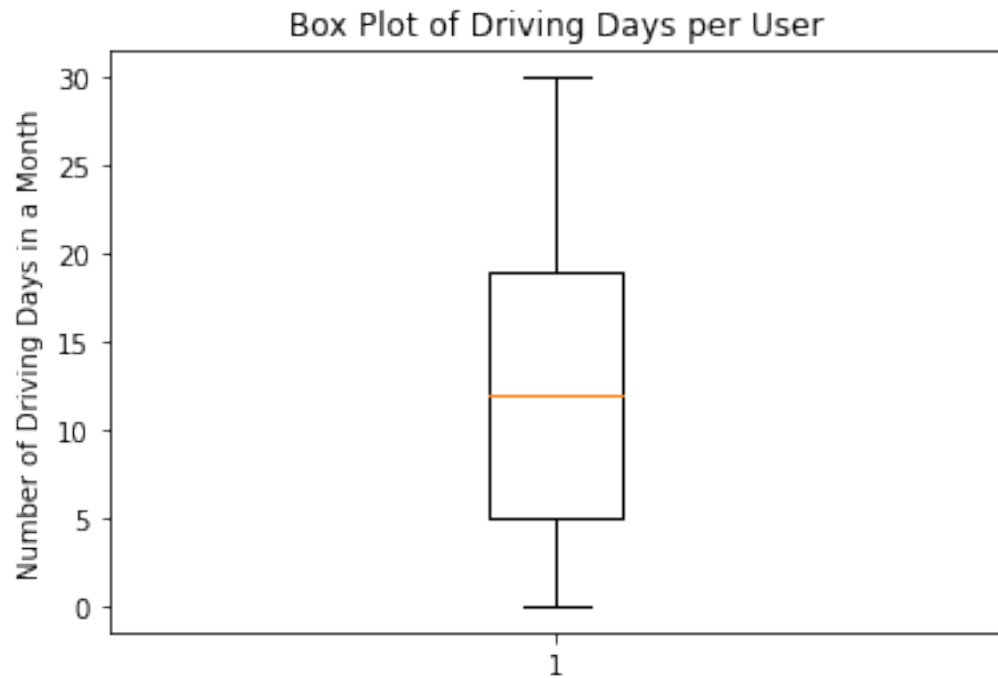


Within the last month, users opened the app a median of 16 times. The box plot reveals a centered distribution. The histogram shows a nearly uniform distribution of ~500 people opening the app on each count of days. However, there are ~250 people who didn't open the app at all and ~250 people who opened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

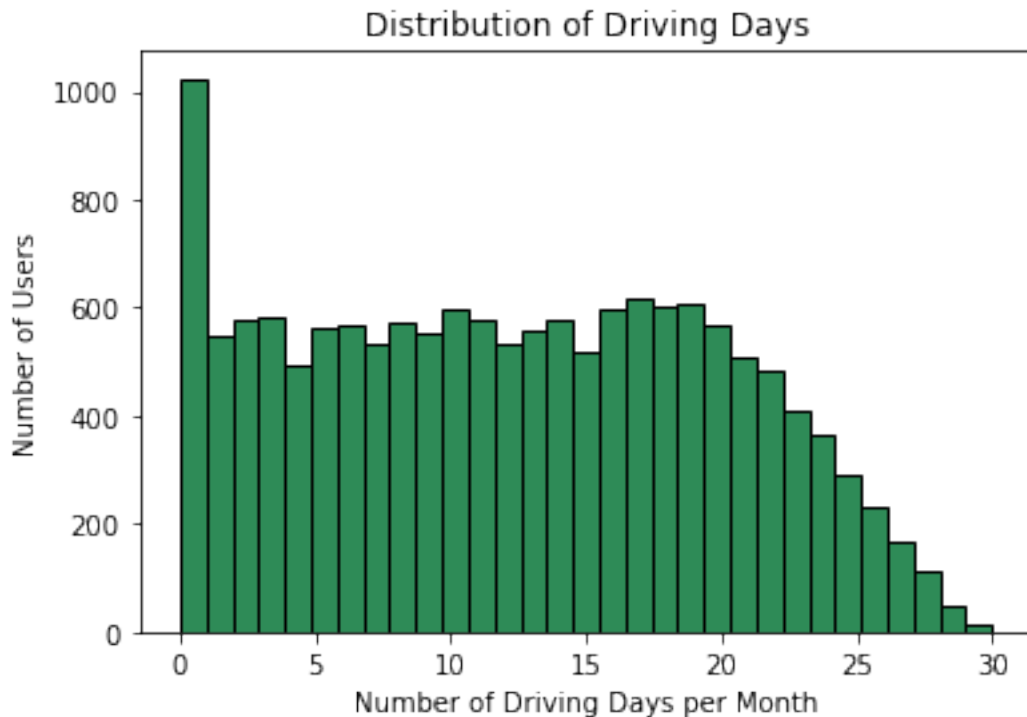
This distribution is noteworthy because it does not mirror the `sessions` distribution, which you might think would be closely correlated with `activity_days`.

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

```
[23]: # Box plot for 'driving_days'
plt.figure(figsize=(6, 4))
plt.boxplot(df['driving_days'])
plt.title('Box Plot of Driving Days per User')
plt.ylabel('Number of Driving Days in a Month')
plt.show()
```



```
[24]: # Histogram for 'driving_days'
plt.figure(figsize=(6, 4))
plt.hist(df['driving_days'], bins=31, color='seagreen', edgecolor='black')
plt.title('Distribution of Driving Days')
plt.xlabel('Number of Driving Days per Month')
plt.ylabel('Number of Users')
plt.show()
```



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.

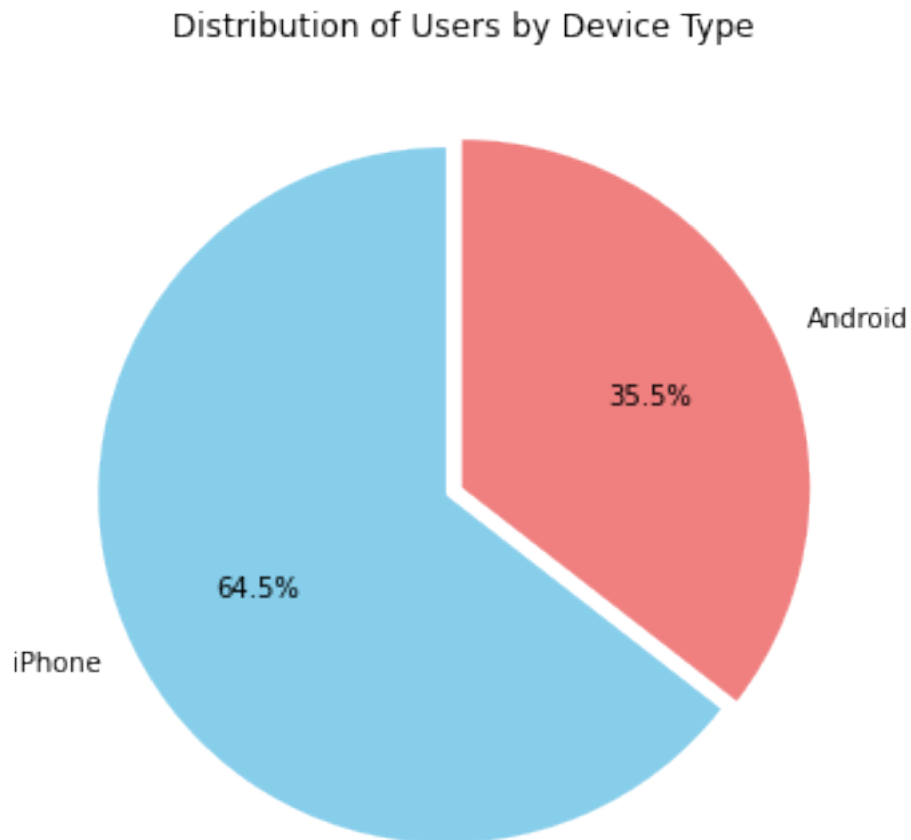
However, there were almost twice as many users (~1,000 vs. ~550) who did not drive at all during the month. This might seem counterintuitive when considered together with the information from `activity_days`. That variable had ~500 users opening the app on each of most of the day counts, but there were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

```
[25]: # Pie chart for 'device'
plt.figure(figsize=(6, 6))
df['device'].value_counts().plot(
    kind='pie',
    autopct='%1.1f%%',
    startangle=90,
    colors=['skyblue', 'lightcoral'],
    explode=(0.05, 0)
)
```

```
plt.title('Distribution of Users by Device Type')
plt.ylabel('') # hides the y-axis label for a cleaner look
plt.show()
```



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*

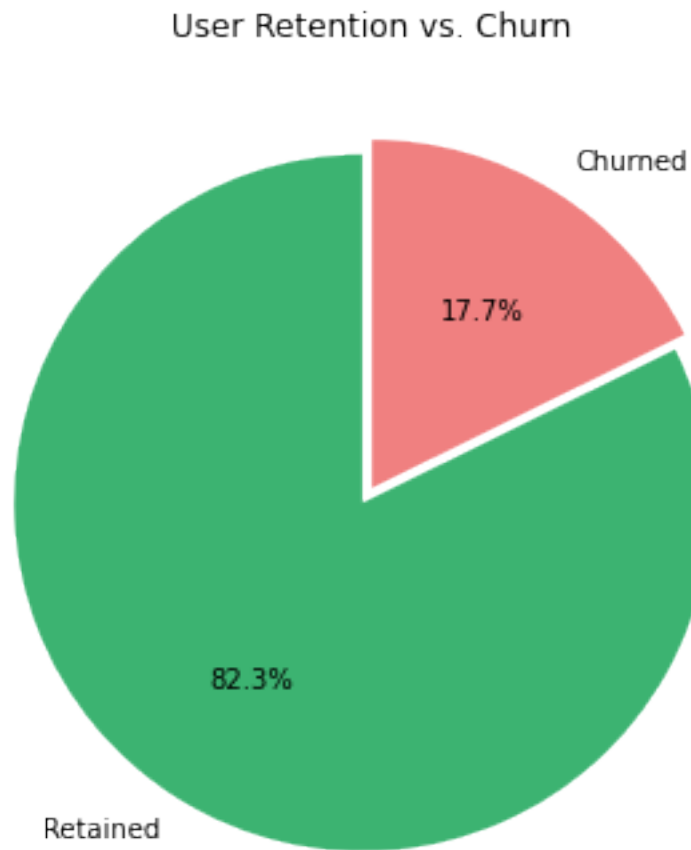
This is also a categorical variable, and as such would not be plotted as a box plot. Plot a pie chart instead.

```
[26]: # Pie chart for 'label'
plt.figure(figsize=(6, 6))
df['label'].value_counts().plot(
    kind='pie',
    autopct='%1.1f%%',
    startangle=90,
```

```

    colors=['mediumseagreen', 'lightcoral'],
    labels=['Retained', 'Churned'],
    explode=(0.05, 0)
)
plt.title('User Retention vs. Churn')
plt.ylabel('') # Hide y-axis label for a cleaner look
plt.show()

```

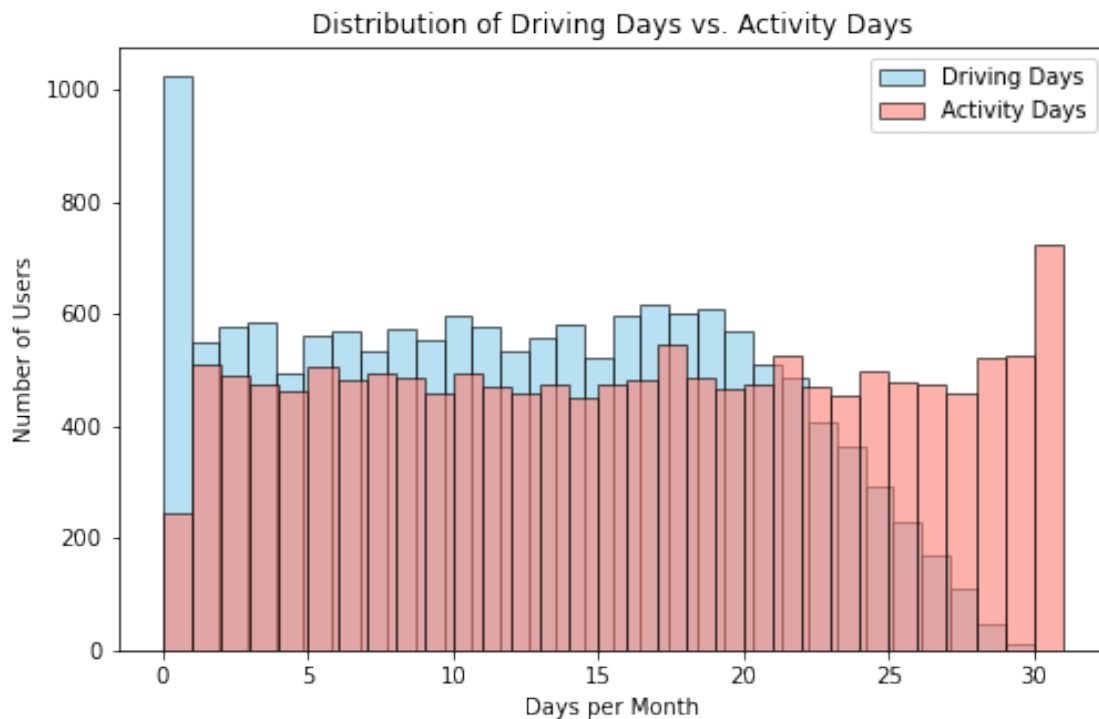


Less than 18% of the users churned.

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`.


```
[27]: # Histogram comparing driving_days and activity_days
plt.figure(figsize=(8, 5))
plt.hist(df['driving_days'], bins=31, alpha=0.6, label='Driving Days',
         color='skyblue', edgecolor='black')
plt.hist(df['activity_days'], bins=31, alpha=0.6, label='Activity Days',
         color='salmon', edgecolor='black')
plt.title('Distribution of Driving Days vs. Activity Days')
plt.xlabel('Days per Month')
plt.ylabel('Number of Users')
plt.legend()
plt.show()
```



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 worthwhile 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`.

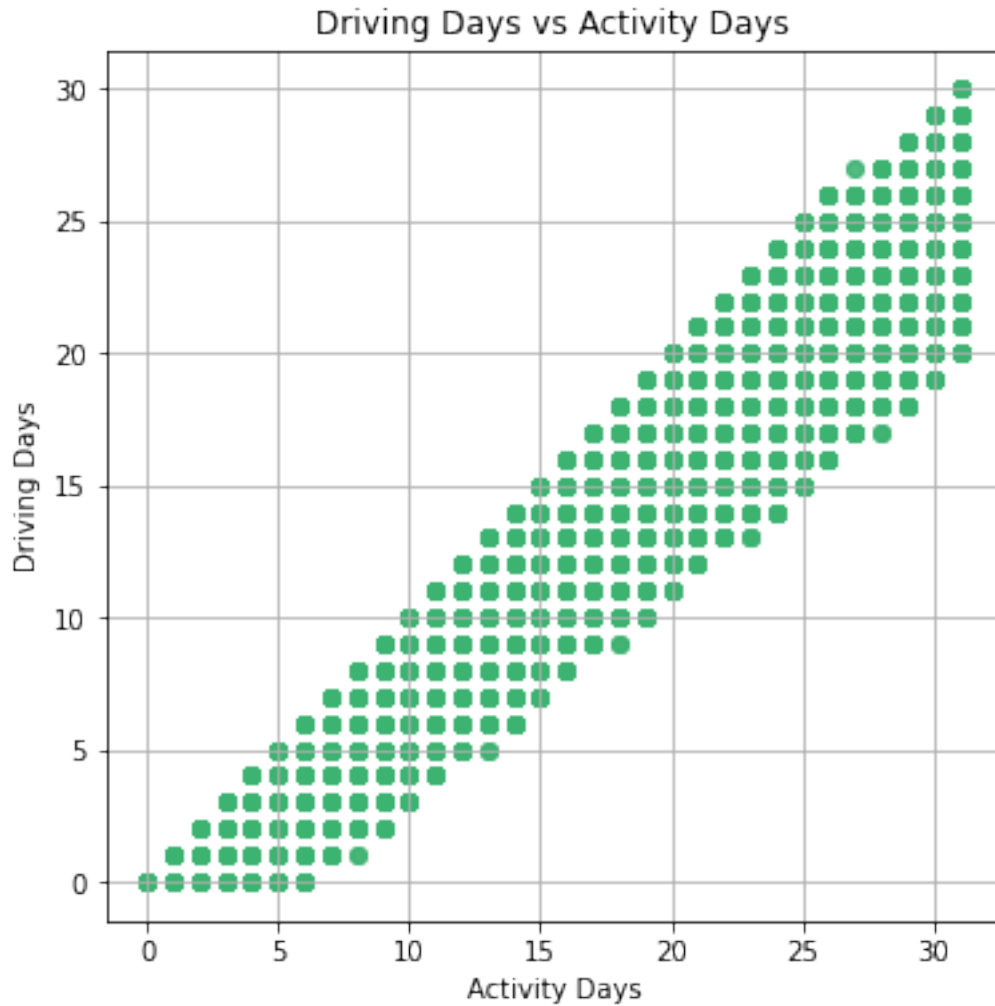
```
[28]: # Confirm the maximum number of days for both variables
print("Max driving_days:", df['driving_days'].max())
print("Max activity_days:", df['activity_days'].max())
```

```
Max driving_days: 30
Max activity_days: 31
```

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.

```
[29]: # Scatter plot of driving_days vs activity_days
plt.figure(figsize=(6, 6))
plt.scatter(df['activity_days'], df['driving_days'], alpha=0.5,
            color='mediumseagreen')
plt.title('Driving Days vs Activity Days')
plt.xlabel('Activity Days')
plt.ylabel('Driving Days')
plt.grid(True)
plt.show()
```



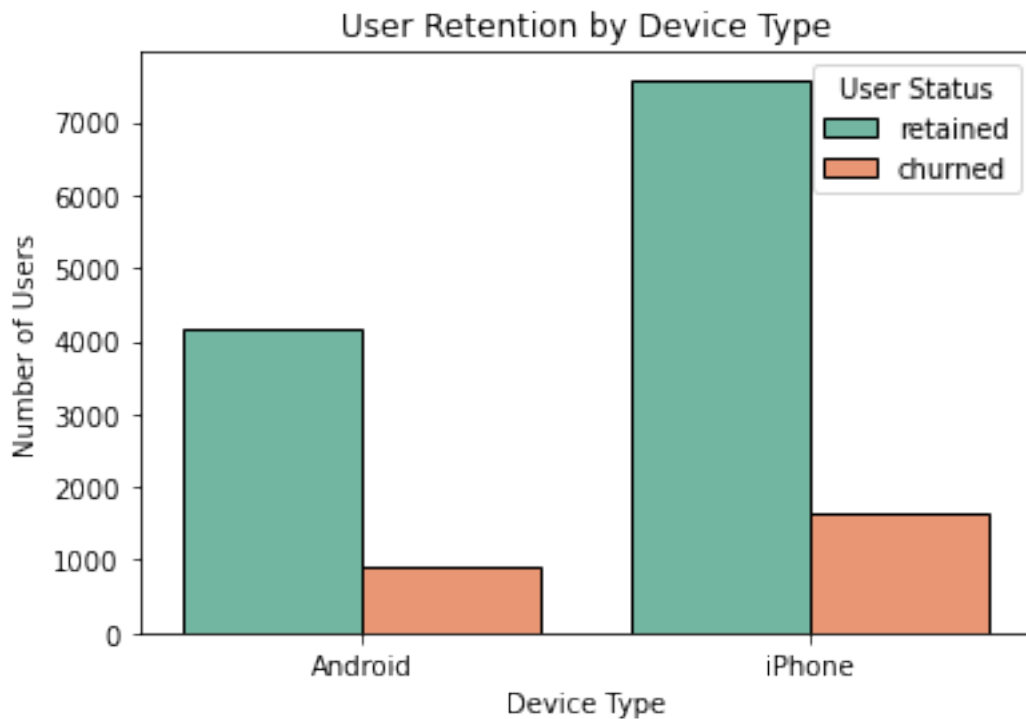
Notice that there is a theoretical limit. If you use the app to drive, then by definition it must count as a day-use as well. In other words, you cannot have more drive-days than activity-days. None of the samples in this data violate this rule, which is good.

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.

```
[31]: # Histogram for retention by device
plt.figure(figsize=(6, 4))

# Create a grouped countplot
import seaborn as sns
sns.countplot(data=df, x='device', hue='label', palette='Set2',
              edgecolor='black')
```

```
plt.title('User Retention by Device Type')
plt.xlabel('Device Type')
plt.ylabel('Number of Users')
plt.legend(title='User Status')
plt.show()
```



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 per driving day last month for users who churned was 697.54 km, versus 289.55 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.

```
[32]: # 1. Create 'km_per_driving_day' column
df['km_per_driving_day'] = df['driven_km_drives'] / df['driving_days']

# 2. Call `describe()` on the new column
df['km_per_driving_day'].describe()
```

```
[32]: count      1.499900e+04
      mean          inf
      std          NaN
      min      3.022063e+00
      25%      1.672804e+02
      50%      3.231459e+02
      75%      7.579257e+02
      max          inf
      Name: km_per_driving_day, dtype: float64
```

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.

```
[33]: # 1. Convert infinite values to zero
      df['km_per_driving_day'].replace(np.inf, 0, inplace=True)

      # 2. Confirm that it worked
      df['km_per_driving_day'].describe()
```

```
[33]: count      14999.000000
      mean       578.963113
      std      1030.094384
      min         0.000000
      25%      136.238895
      50%      272.889272
      75%      558.686918
      max     15420.234110
      Name: km_per_driving_day, dtype: float64
```

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`.

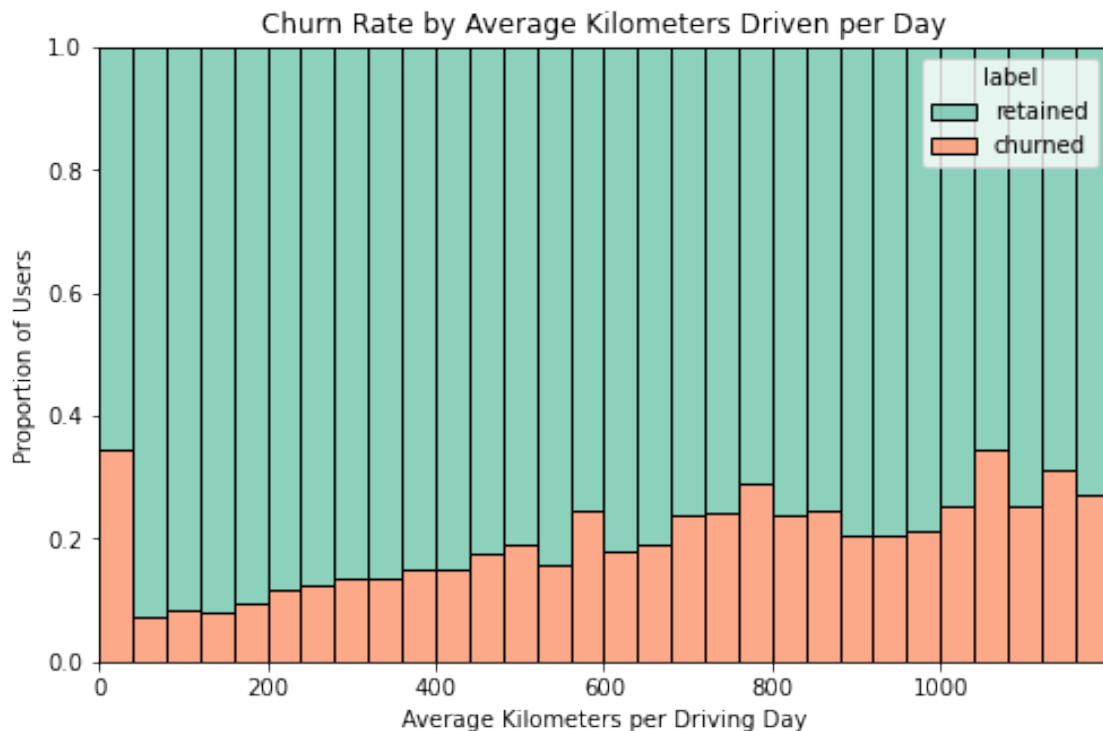
```
[34]: # Histogram for km_per_driving_day (churn vs retained)
      plt.figure(figsize=(8, 5))
      sns.histplot(
          data=df[df['km_per_driving_day'] <= 1200], # exclude unrealistic values
```

```

x='km_per_driving_day',
hue='label',
multiple='fill',
bins=30,
palette='Set2'
)

plt.title('Churn Rate by Average Kilometers Driven per Day')
plt.xlabel('Average Kilometers per Driving Day')
plt.ylabel('Proportion of Users')
plt.show()

```



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.

```

[35]: # Histogram for churn rate by number of driving days
plt.figure(figsize=(8, 5))
sns.histplot(
    data=df,

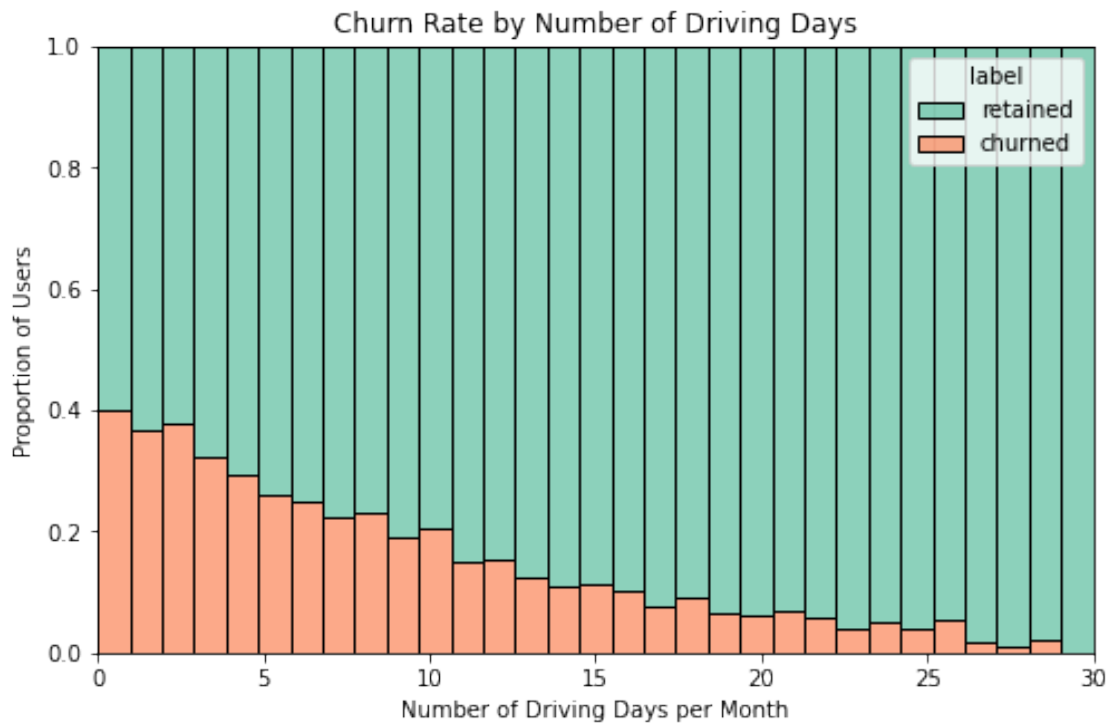
```

```

x='driving_days',
hue='label',
multiple='fill',
bins=31,
palette='Set2'
)

plt.title('Churn Rate by Number of Driving Days')
plt.xlabel('Number of Driving Days per Month')
plt.ylabel('Proportion of Users')
plt.show()

```



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.

```
[36]: # 1. Create 'percent_sessions_in_last_month' column
df['percent_sessions_in_last_month'] = (df['sessions'] / df['total_sessions']) * 100

# 2. Check basic statistics
df['percent_sessions_in_last_month'].describe()
```

```
[36]: count      14999.000000
      mean         44.925534
      std         28.691863
      min          0.000000
      25%         19.622145
      50%         42.309703
      75%         68.721626
      max         153.063707
      Name: percent_sessions_in_last_month, dtype: float64
```

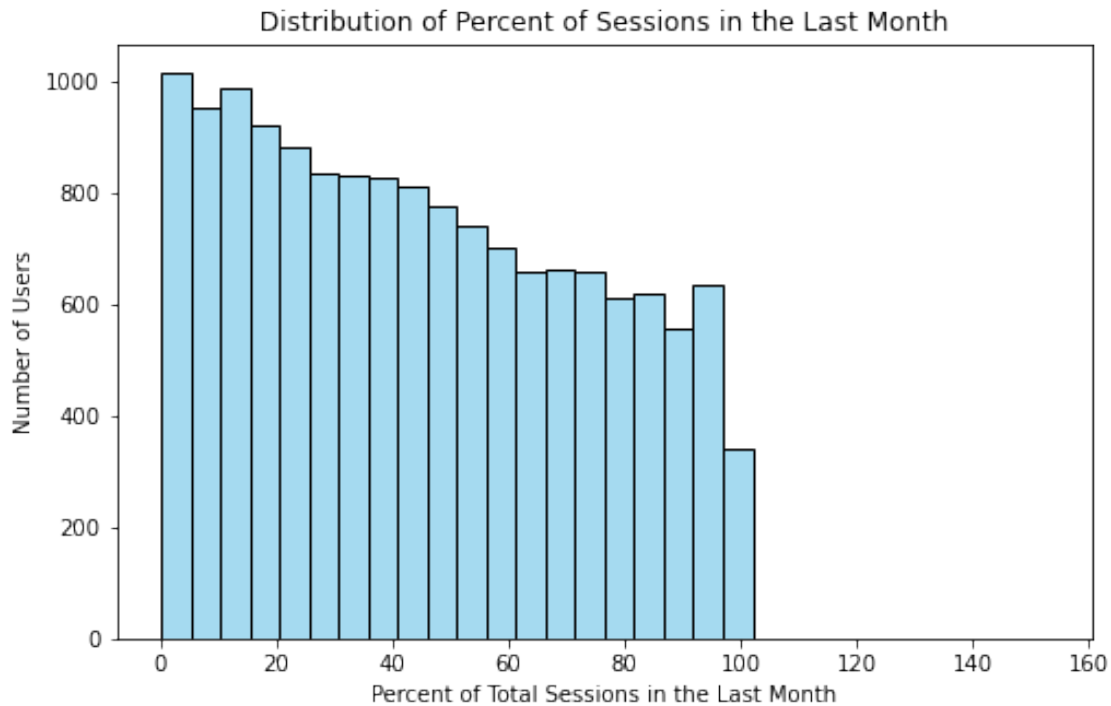
What is the median value of the new column?

```
[37]: # Median value
median_percent = df['percent_sessions_in_last_month'].median()
print("Median percent of sessions in the last month:", median_percent)
```

Median percent of sessions in the last month: 42.30970299276318

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

```
[38]: # Histogram of percent_sessions_in_last_month
plt.figure(figsize=(8, 5))
sns.histplot(df['percent_sessions_in_last_month'], bins=30, kde=False, color='skyblue', edgecolor='black')
plt.title('Distribution of Percent of Sessions in the Last Month')
plt.xlabel('Percent of Total Sessions in the Last Month')
plt.ylabel('Number of Users')
plt.show()
```

Check the median value of the `n_days_after_onboarding` variable.

```
[39]: # Median days since onboarding
median_days = df['n_days_after_onboarding'].median()
print("Median number of days since onboarding:", median_days)
```

Median number of days since onboarding: 1741.0

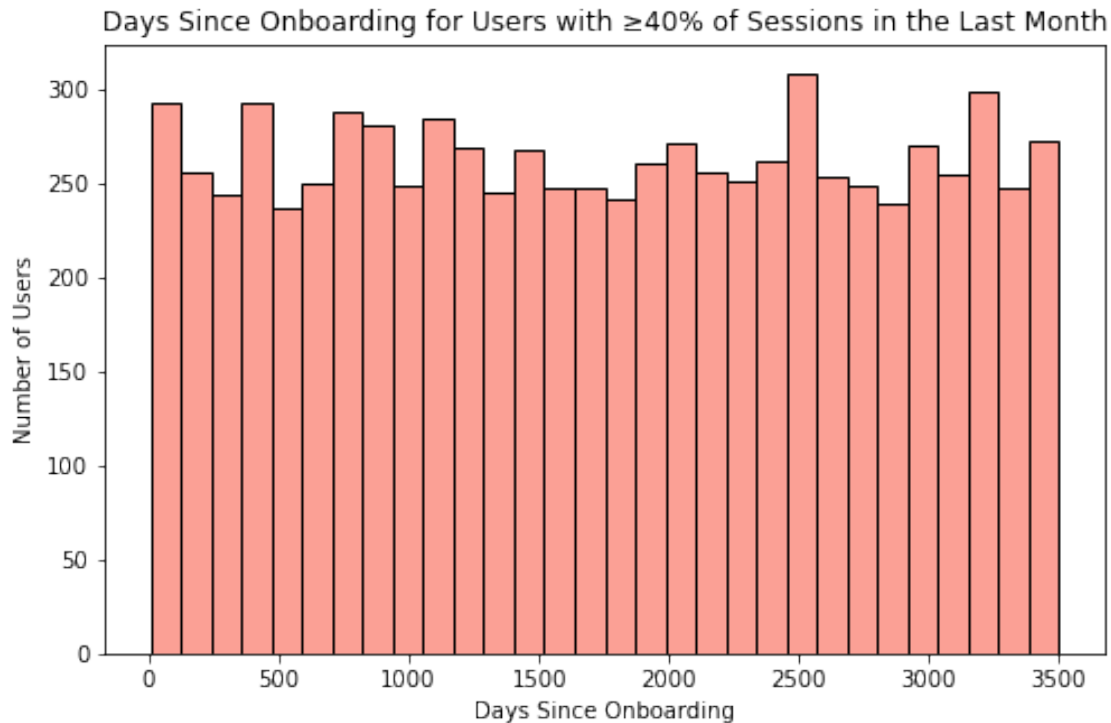
Half of the people in the dataset had 40% or more of their sessions in just the last month, yet the overall median time since onboarding is almost five years.

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.

```
[40]: # Filter users who had 40% or more of their sessions in the last month
recently_active = df[df['percent_sessions_in_last_month'] >= 40]

# Histogram of n_days_after_onboarding for these users
plt.figure(figsize=(8, 5))
sns.histplot(recently_active['n_days_after_onboarding'], bins=30,
             color='salmon', edgecolor='black')
plt.title('Days Since Onboarding for Users with 40% of Sessions in the Last Month')
plt.xlabel('Days Since Onboarding')
plt.ylabel('Number of Users')
```

```
plt.show()
```



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.

```
[41]: # Function to cap outliers at the 95th percentile
def cap_outliers(series):
    threshold = series.quantile(0.95)
    series = np.where(series > threshold, threshold, series)
```

```
return series
```

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

```
[42]: # Apply the outlier capping function
cols_to_cap = ['sessions', 'drives', 'total_sessions', 'driven_km_drives',
               'duration_minutes_drives']

for col in cols_to_cap:
    df[col] = cap_outliers(df[col])
```

Call describe() to see if your change worked.

```
[43]: # Confirm that outliers were capped
df[cols_to_cap].describe()
```

```
[43]:
```

	sessions	drives	total_sessions	driven_km_drives	\
count	14999.000000	14999.000000	14999.000000	14999.000000	
mean	76.568705	64.058204	184.031320	3939.632764	
std	67.297958	55.306924	118.600463	2216.041510	
min	0.000000	0.000000	0.220211	60.441250	
25%	23.000000	20.000000	90.661156	2212.600607	
50%	56.000000	48.000000	159.568115	3493.858085	
75%	112.000000	93.000000	254.192341	5289.861262	
max	243.000000	201.000000	454.363204	8889.794236	

	duration_minutes_drives
count	14999.000000
mean	1789.647426
std	1222.705167
min	18.282082
25%	835.996260
50%	1478.249859
75%	2464.362632
max	4668.899349

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?

I have learned... Through this exploratory data analysis, I learned that the overall churn rate among Waze users is approximately 17%, and that this rate is consistent between iPhone and Android users. Users who drive longer distances per driving day tend to be more likely to churn, while those who drive more frequently (more driving days) are less likely to churn.

Additionally, nearly half of the users completed over 40% of their total sessions in the most recent month, suggesting a sudden spike in engagement — even among long-term users who have been active for several years.

My other questions are...

Why did so many long-term users suddenly increase their activity during the last month?

What factors (e.g., new app features, seasonality, or marketing campaigns) might explain this surge in usage?

What specific factors drive churn among long-distance drivers — is it related to battery drain, navigation accuracy, or trip type?

Are there regional or temporal patterns in churn (e.g., certain countries or months)?

My client would likely want to know...

How can Waze reduce churn among long-distance drivers?

What strategies could sustain engagement among recently reactivated long-term users?

Whether usage patterns (driving days vs. driving distance) can be used to predict future churn and enable targeted retention campaigns.

How to translate these findings into actionable product changes or marketing efforts, such as personalized notifications or improved routing for heavy users.

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.

```
[44]: # Correlation heatmap across numeric features + churn as 0/1
import numpy as np, pandas as pd, seaborn as sns, matplotlib.pyplot as plt
```

```

# Encode churn: 1 = churned, 0 = retained (temp column for analysis)
df['_label_num'] = (df['label'].str.lower() == 'churned').astype(int)

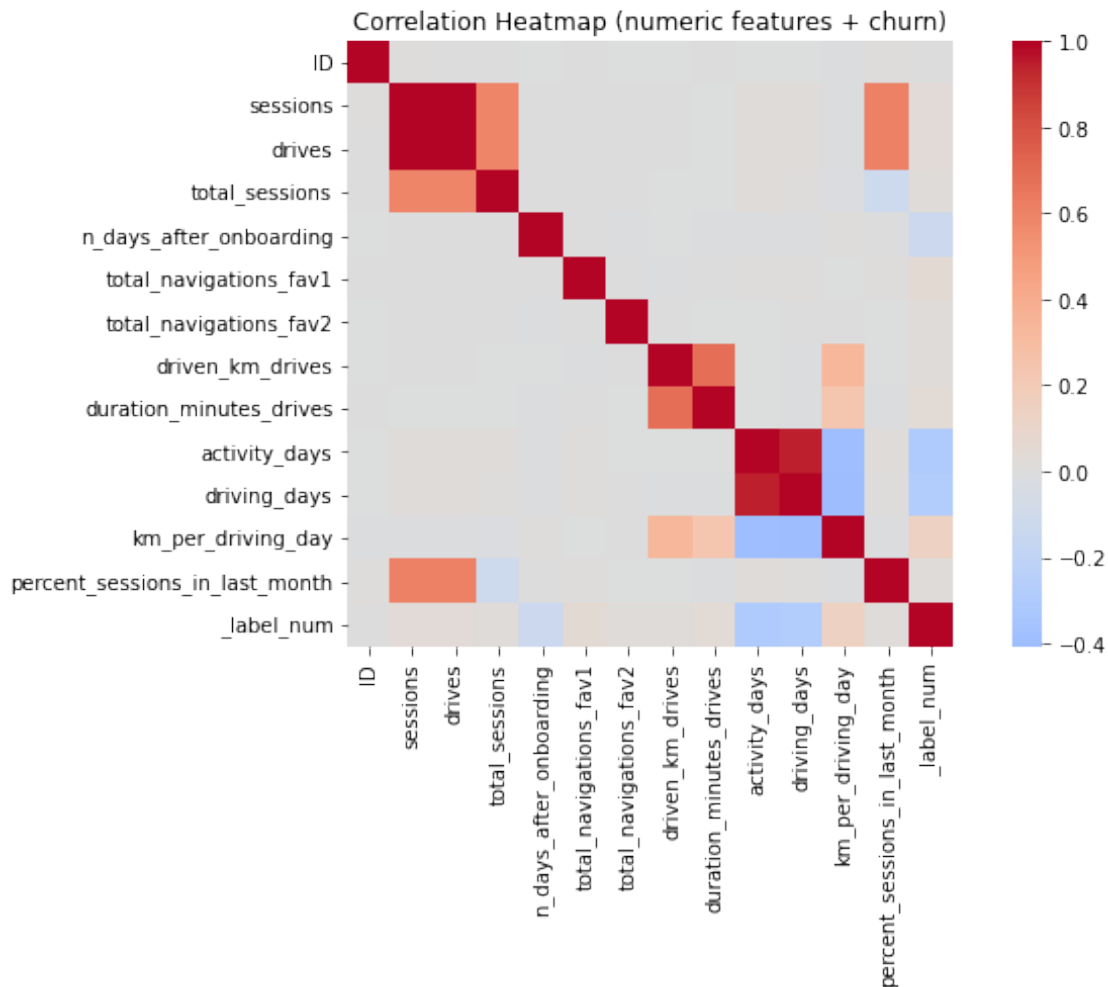
# Pick numeric columns automatically, add churn numeric
num_cols = df.select_dtypes(include=['number']).columns.tolist()
if '_label_num' not in num_cols:
    num_cols +=['_label_num']

corr = df[num_cols].corr()

plt.figure(figsize=(10,7))
sns.heatmap(corr, annot=False, cmap='coolwarm', center=0, square=True)
plt.title('Correlation Heatmap (numeric features + churn)')
plt.tight_layout()
plt.show()

# Optional: print top correlations with churn
print(
    corr['_label_num']
    .drop('_label_num')
    .sort_values(ascending=False)
    .to_frame('corr_with_churn')
    .head(10)
)

```



	corr_with_churn
km_per_driving_day	0.146798
total_navigations_fav1	0.054237
duration_minutes_drives	0.040527
drives	0.034528
sessions	0.033695
total_sessions	0.022467
driven_km_drives	0.020131
total_navigations_fav2	0.019121
percent_sessions_in_last_month	0.018415
ID	0.004777

```
[45]: # Churn rate by driving frequency buckets + boxplot to compare distance-per-day
      ↪ by churn
      import matplotlib.pyplot as plt
      import seaborn as sns
```

```

import numpy as np
import pandas as pd

# Make bins for driving_days
bins = [-0.1, 0, 5, 10, 20, 25, 30, 31]
labels = ['0', '1-5', '6-10', '11-20', '21-25', '26-30', '31']
tmp = df.copy()
tmp['driving_days_bin'] = pd.cut(tmp['driving_days'], bins=bins, labels=labels,
    ↪right=True)

# Compute churn rate per bin
rate = (tmp.groupby('driving_days_bin')['label']
        .apply(lambda s: (s.str.lower()=='churned').mean())
        .rename('churn_rate')
        .reset_index())

# Ensure km_per_driving_day exists (recreate if needed) and clean impossible_
    ↪values
if 'km_per_driving_day' not in tmp.columns:
    tmp['km_per_driving_day'] = tmp['driven_km_drives'] / tmp['driving_days']
tmp['km_per_driving_day'].replace([np.inf, -np.inf], 0, inplace=True)
tmp = tmp[tmp['km_per_driving_day'] <= 1200]

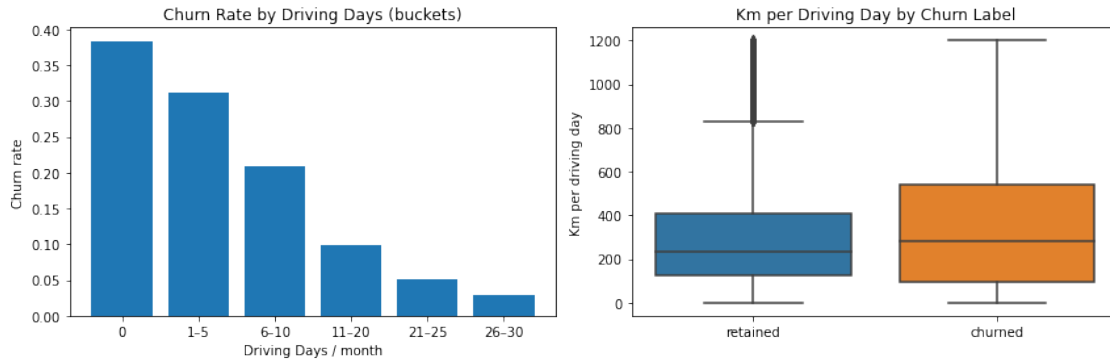
# Plot
fig, ax = plt.subplots(1, 2, figsize=(12,4))

# (A) Churn rate by driving_days bucket
ax[0].bar(rate['driving_days_bin'].astype(str), rate['churn_rate'])
ax[0].set_title('Churn Rate by Driving Days (buckets)')
ax[0].set_xlabel('Driving Days / month')
ax[0].set_ylabel('Churn rate')

# (B) Distance-per-driving-day by churn label
sns.boxplot(data=tmp, x='label', y='km_per_driving_day', ax=ax[1])
ax[1].set_title('Km per Driving Day by Churn Label')
ax[1].set_xlabel('')
ax[1].set_ylabel('Km per driving day')

plt.tight_layout()
plt.show()

```



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?
1. Most variables like sessions and drives were right-skewed, meaning a few users were highly active while most showed moderate activity. `n_days_after_onboarding` was fairly uniform, showing users with a wide range of tenures.
 2. No major data errors were found, but outliers existed (e.g., users driving unrealistically long distances). These were handled by capping at the 95th percentile.
 3. Further questions include why long-term users increased their activity recently, and why long-distance drivers churned more often. External factors such as seasonality or app changes may explain this.
 4. About 17% of users churned, while 83% were retained, with similar rates for Android and iPhone users.

5. Higher driving distance per day correlated with higher churn, while more driving days and sessions correlated with lower churn — indicating frequent app use promotes retention.
6. New and long-term users were both well represented. Interestingly, many long-term users logged a large share of sessions recently, showing renewed engagement.

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.