

Activity_Course 3 Automatidata project lab

October 22, 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 Rodriguez, 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&uniquifier=1)
- [Analyze] (#scrollTo=mA7Mz_SnI8km&line=4&uniquifier=1)
- [Construct] (#scrollTo=Lca9c8XON8lc&line=2&uniquifier=1)
- [Execute] (#scrollTo=401PgchTPr4E&line=2&uniquifier=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.

```
[14]: # Import packages and libraries
      #==> ENTER YOUR CODE HERE
      import pandas as pd
      import matplotlib.pyplot as plt
      import numpy as np
      import datetime as dt
      import seaborn as sns
```

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.

```
[15]: # 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.

```
[3]: 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

```
[4]: df.size
```

```
[4]: 408582
```

Use describe...

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

```
[12]:
```

	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

And info.

```
[13]: 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
```

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

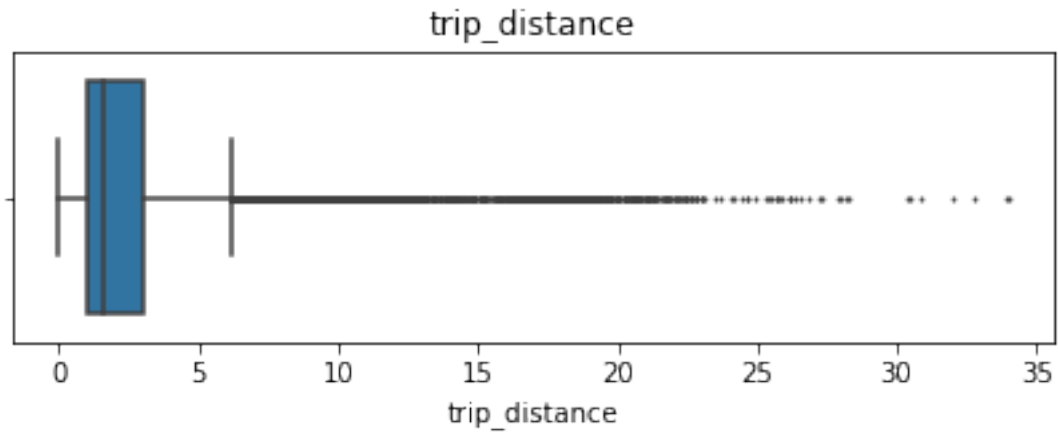
Perform a check for outliers on relevant columns such as trip distance and trip duration. Remember, some of the best ways to identify the presence of outliers in data are box plots and histograms.

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

```
[16]: df['tpep_pickup_datetime']=pd.to_datetime(df['tpep_pickup_datetime'])  
      df['tpep_dropoff_datetime']=pd.to_datetime(df['tpep_dropoff_datetime'])
```

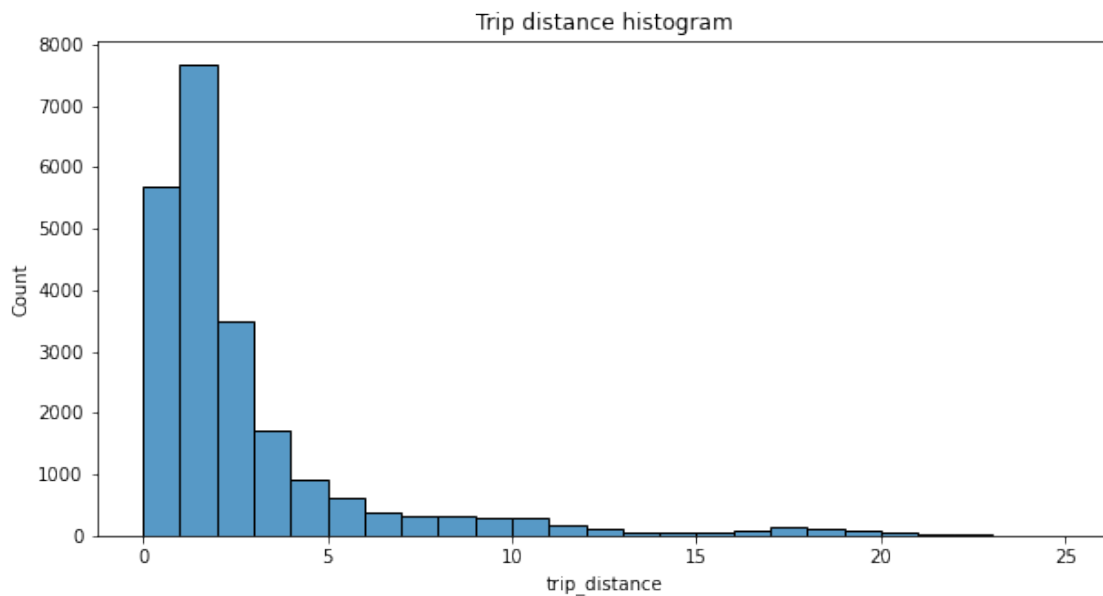
trip distance

```
[5]: plt.figure(figsize=(7,2))  
      plt.title('trip_distance')  
      sns.boxplot(data=None, x=df['trip_distance'], fliersize=1);
```



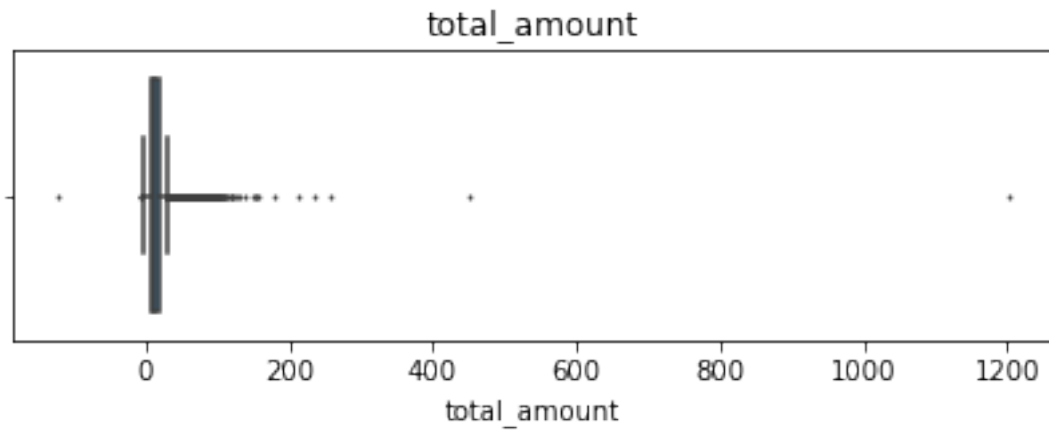
```
[6]: plt.figure(figsize=(10,5))
plt.title('Trip distance histogram');
sns.histplot(df['trip_distance'], bins=range(0,26,1))
```

```
[6]: <matplotlib.axes._subplots.AxesSubplot at 0x7920b5ee2450>
```



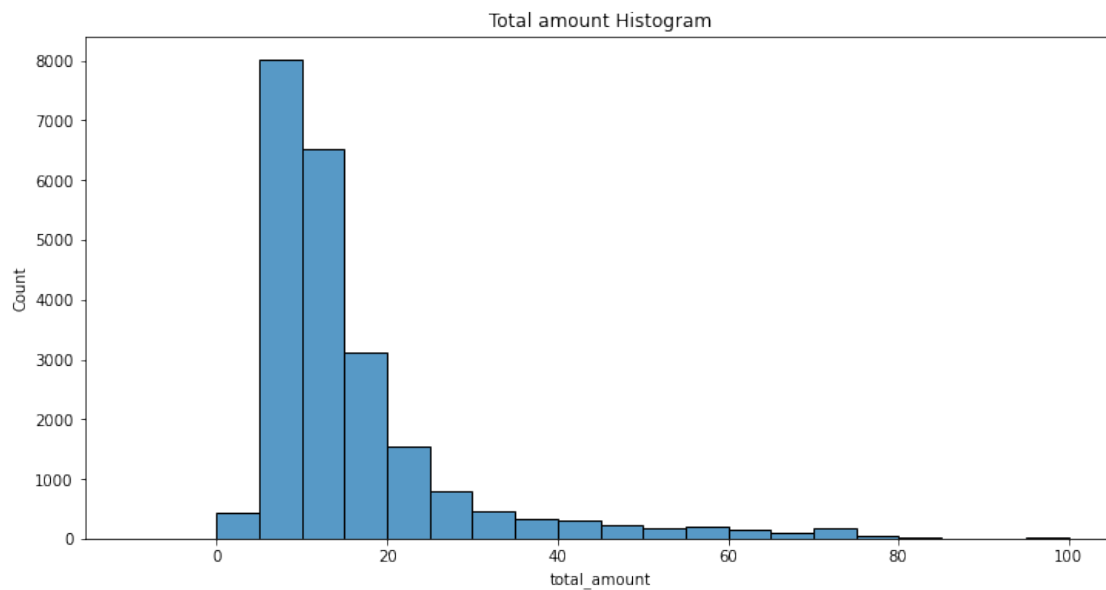
total amount

```
[9]: plt.figure(figsize=(7,2))
plt.title('total_amount')
sns.boxplot(data=None, x=df['total_amount'], fliersize=1);
```

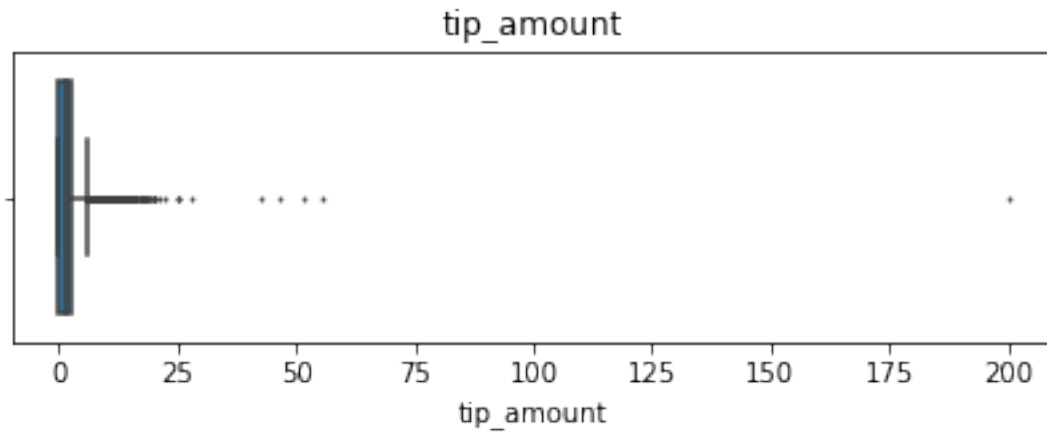
```
[10]: plt.figure(figsize=(12,6))
plt.title('Total amount Histogram')
sns.histplot(df['total_amount'], bins=range(-10,101,5))
```

```
[10]: <matplotlib.axes._subplots.AxesSubplot at 0x7920b67df8d0>
```

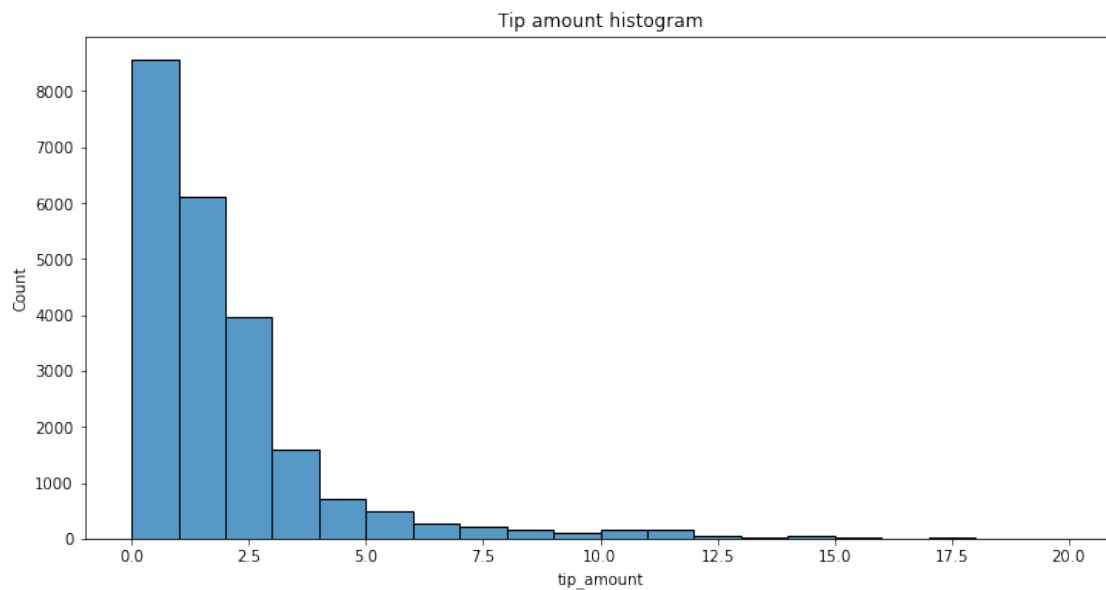


tip amount

```
[11]: plt.figure(figsize=(7,2))
plt.title('tip_amount')
sns.boxplot(x=df['tip_amount'], fliersize=1);
```



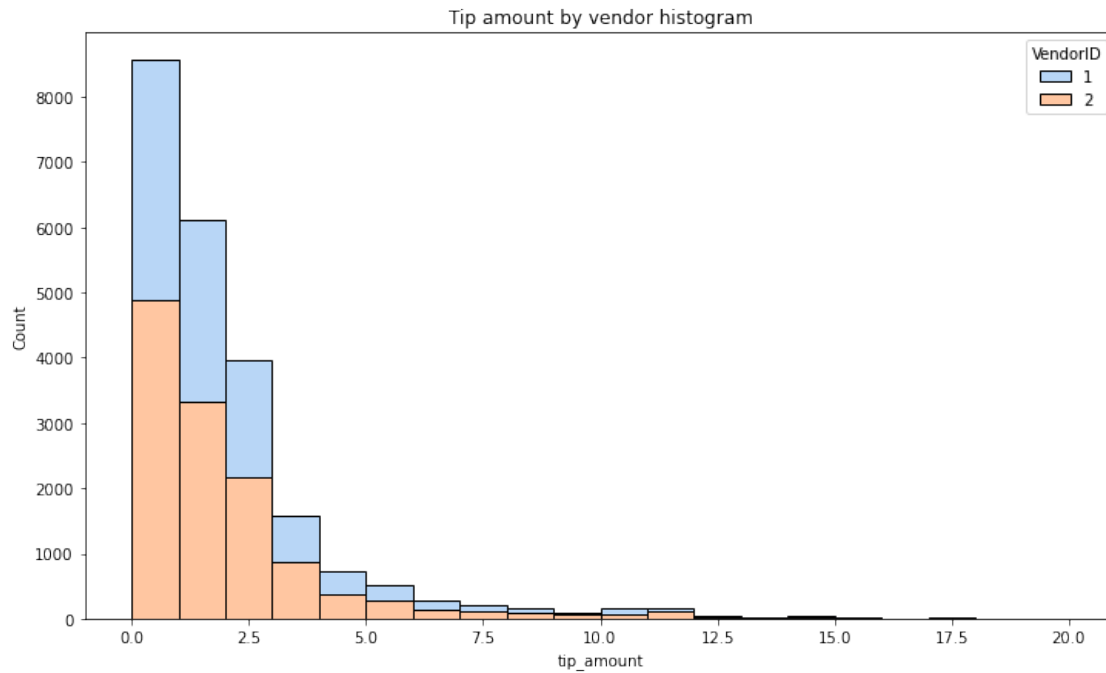
```
[13]: plt.figure(figsize=(12,6))
sns.histplot(df['tip_amount'], bins=range(0,21,1))
plt.title('Tip amount histogram');
```



tip_amount by vendor

```
[15]: plt.figure(figsize=(12,7))
sns.histplot(data=df, x='tip_amount', bins=range(0,21,1),
             hue='VendorID',
             multiple='stack',
             palette='pastel')
```

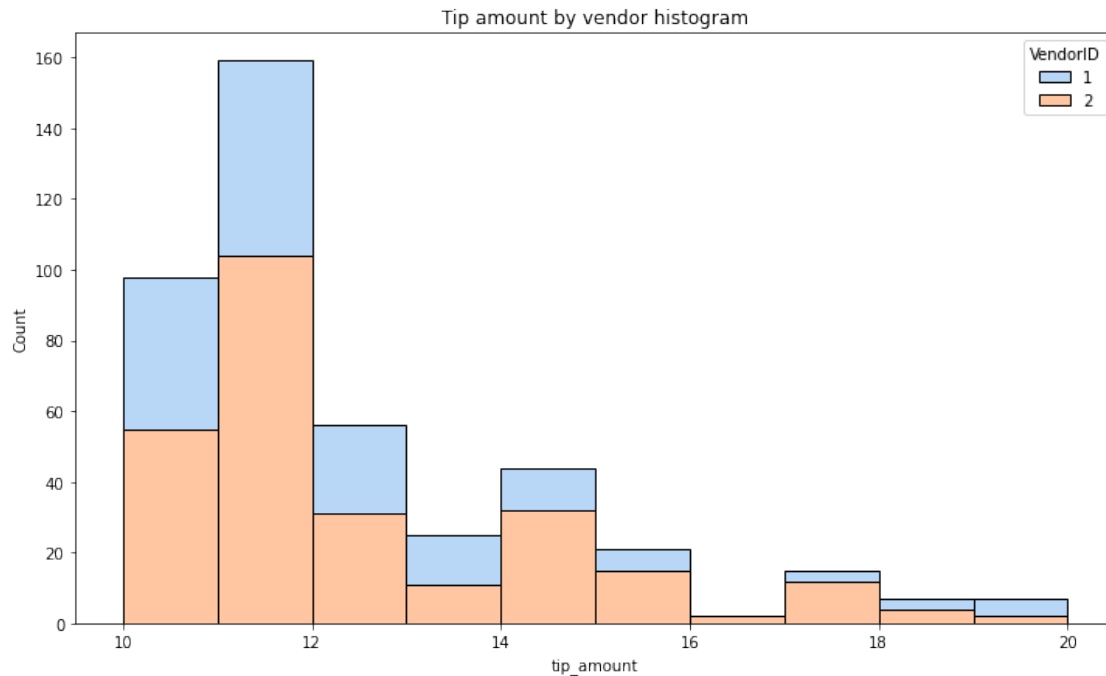
```
plt.title('Tip amount by vendor histogram');
```



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.

```
[11]: tips_over_ten = df[df['tip_amount'] > 10]
plt.figure(figsize=(12,7))
sns.histplot(data=tips_over_ten, x='tip_amount', bins=range(10,21,1),
             hue='VendorID',
             multiple='stack',
             palette='pastel')

plt.title('Tip amount by vendor histogram');
```



Mean tips by passenger count

Examine the unique values in the `passenger_count` column.

```
[5]: df['passenger_count'].value_counts()
```

```
[5]: 1    16117
      2     3305
      5     1143
      3     953
      6     693
      4     455
      0       33
      Name: passenger_count, dtype: int64
```

```
[6]: # Calculatin mean tips by passenger_count
mean_tips_by_passenger_count = df.groupby(['passenger_count']).
    ↳mean()['tip_amount']
mean_tips_by_passenger_count
```

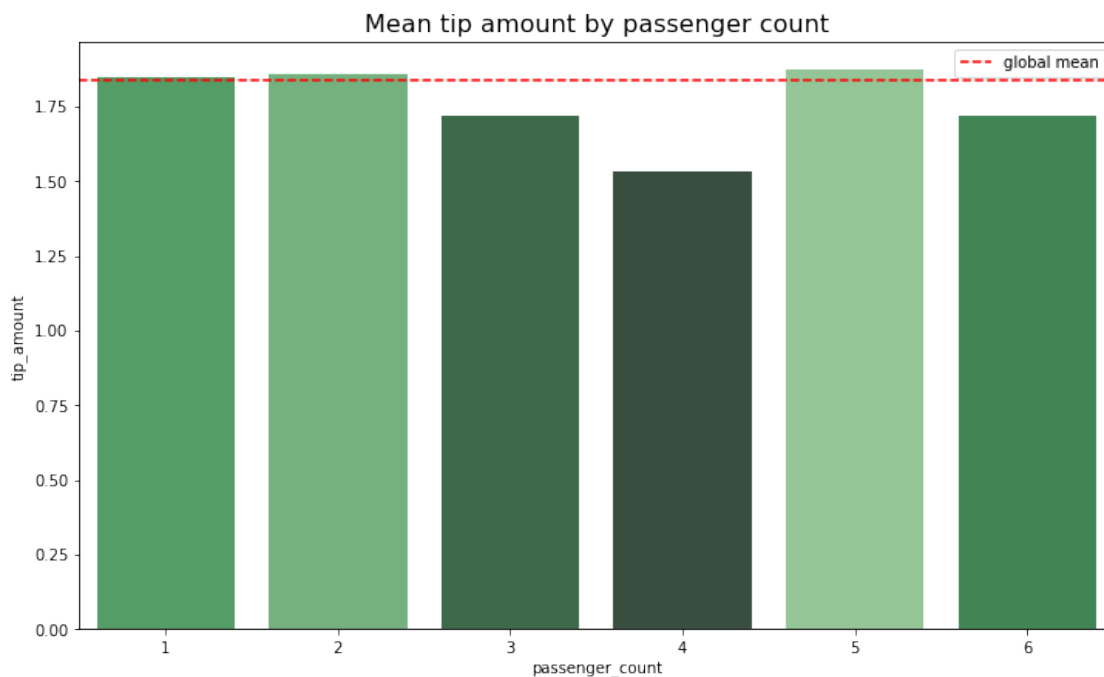
```
[6]:          tip_amount
passenger_count
0          2.135758
1          1.848920
2          1.856378
3          1.716768
```

4	1.530264
5	1.873185
6	1.720260

```
[10]: # bar plot for mean tips by passenger count
data = mean_tips_by_passenger_count.tail(-1)
pal = sns.color_palette("Greens_d", len(data))
rank = data['tip_amount'].argsort().argsort()

plt.figure(figsize=(12,7))
plt.title('Mean tip amount by passenger count', fontsize=16);
ax = sns.barplot(x=data.index,
                 y=data['tip_amount'],
                 palette=np.array(pal[:,-1])[rank])
ax.axhline(df['tip_amount'].mean(), ls='--', color='red', label='global mean')
ax.legend()
```

[10]: <matplotlib.legend.Legend at 0x7c90df304b90>



Create month and day columns

```
[17]: df['month'] = df['tpep_pickup_datetime'].dt.month_name()

df['day'] = df['tpep_pickup_datetime'].dt.day_name()
```

Plot total ride count by month

Begin by calculating total ride count by month.

```
[18]: monthly_rides = df['month'].value_counts()  
monthly_rides
```

```
[18]: March          2049  
      October       2027  
      April         2019  
      May           2013  
      January       1997  
      June          1964  
      December      1863  
      November      1843  
      February      1769  
      September     1734  
      August        1724  
      July          1697  
      Name: month, dtype: int64
```

Reorder the results to put the months in calendar order.

```
[19]: month_order = ['January', 'February', 'March', 'April', 'May', 'June', 'July',  
                    'August', 'September', 'October', 'November', 'December']  
  
monthly_rides = monthly_rides.reindex(index=month_order)  
monthly_rides
```

```
[19]: January         1997  
      February       1769  
      March          2049  
      April          2019  
      May            2013  
      June           1964  
      July           1697  
      August         1724  
      September     1734  
      October        2027  
      November      1843  
      December      1863  
      Name: month, dtype: int64
```

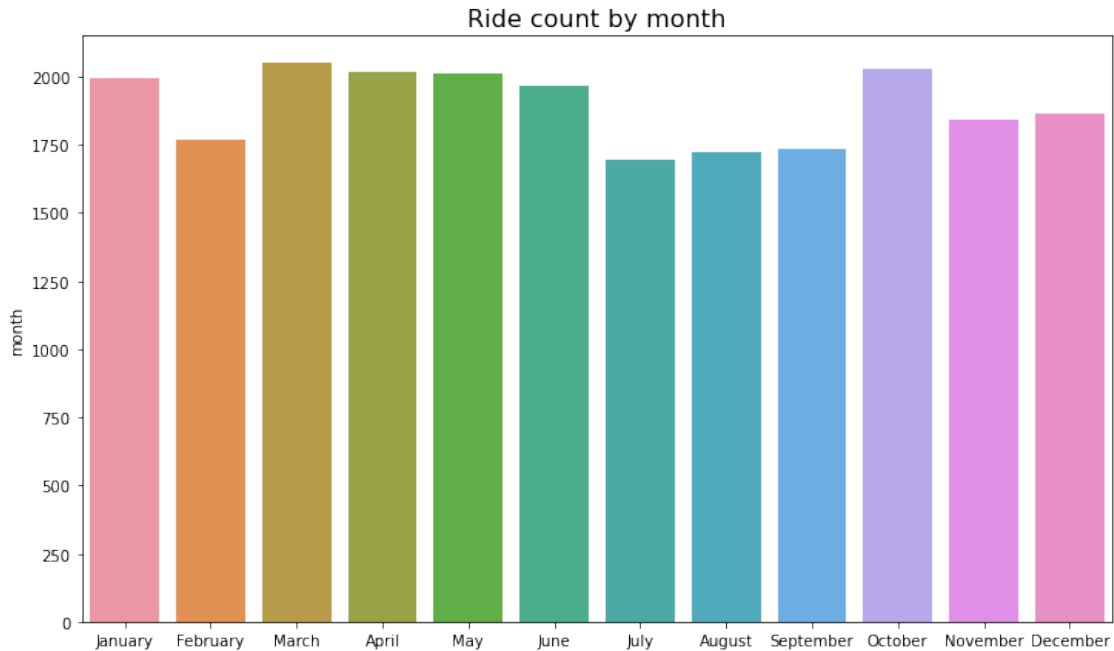
```
[20]: monthly_rides.index
```

```
[20]: Index(['January', 'February', 'March', 'April', 'May', 'June', 'July',  
        'August', 'September', 'October', 'November', 'December'],  
       dtype='object')
```

```
[21]: plt.figure(figsize=(12,7))
plt.title('Ride count by month', fontsize=16);

sns.barplot(x=monthly_rides.index, y=monthly_rides)
```

```
[21]: <matplotlib.axes._subplots.AxesSubplot at 0x7c90df03f2d0>
```



Plot total ride count by day

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

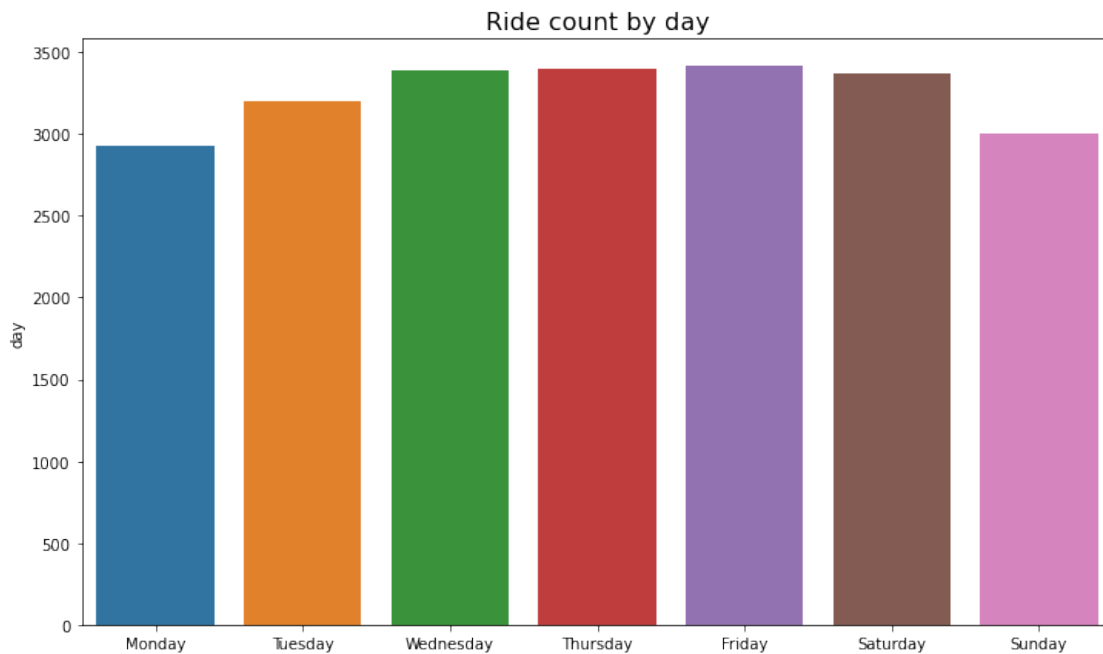
```
[23]: daily_rides = df['day'].value_counts()
day_order = ['Monday', 'Tuesday', 'Wednesday', 'Thursday', 'Friday', 'Saturday', 'Sunday']
daily_rides = daily_rides.reindex(index=day_order)
daily_rides
```

```
[23]: Monday      2931
      Tuesday    3198
      Wednesday  3390
      Thursday   3402
      Friday     3413
      Saturday   3367
      Sunday     2998
      Name: day, dtype: int64
```

```
[24]: # barplot for daily_rides

plt.figure(figsize=(12,7))
plt.title('Ride count by day', fontsize=16);
sns.barplot(x=daily_rides.index, y=daily_rides)
```

[24]: <matplotlib.axes._subplots.AxesSubplot at 0x7c90def56bd0>



Plot total revenue by day of the week

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

```
[25]: total_revenue_day = df.groupby('day').sum()[['total_amount']]

day_order = ['Monday', 'Tuesday', 'Wednesday', 'Thursday', 'Friday',
             ↪ 'Saturday', 'Sunday']

total_revenue_day = total_revenue_day.reindex(index=day_order)
total_revenue_day
```

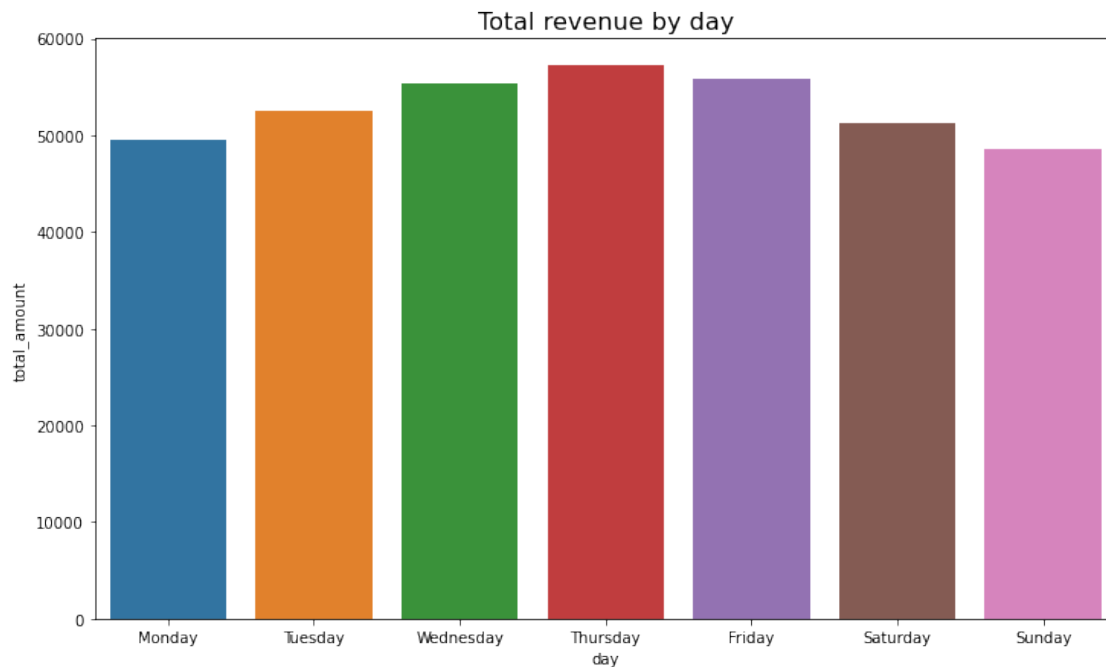
```
[25]:
```

day	total_amount
Monday	49574.37
Tuesday	52527.14
Wednesday	55310.47
Thursday	57181.91

Friday	55818.74
Saturday	51195.40
Sunday	48624.06

```
[27]: plt.figure(figsize=(12,7))
sns.barplot(x=total_revenue_day.index, y=total_revenue_day['total_amount'])

plt.title('Total revenue by day', fontsize=16);
```



Plot total revenue by month

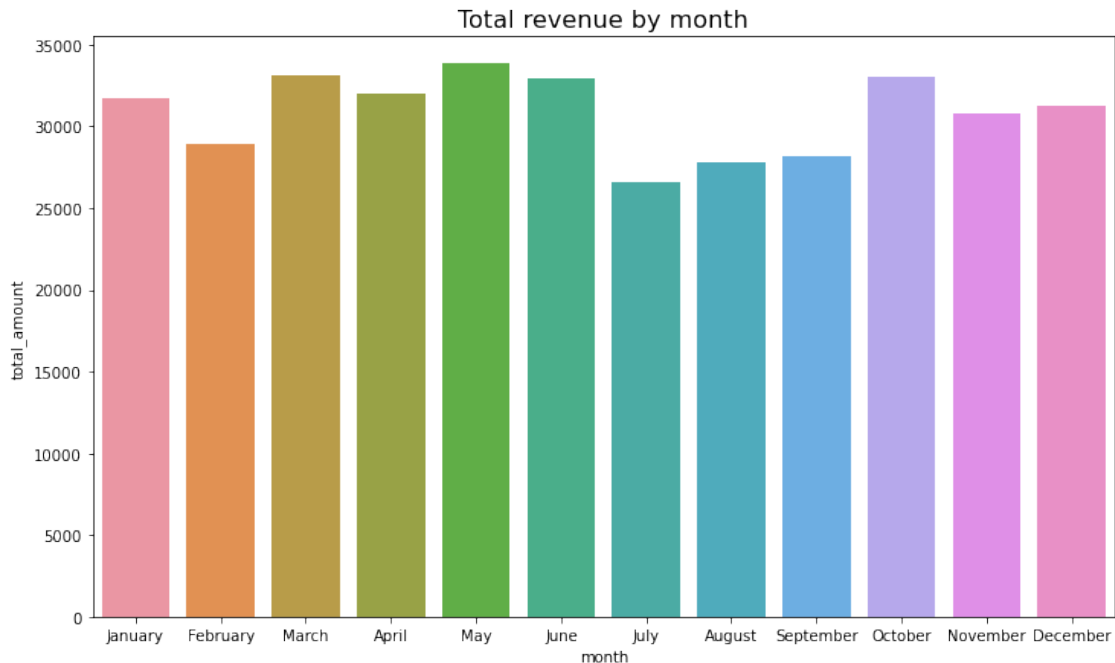
```
[28]: total_amount_month = df.groupby('month').sum()[['total_amount']]
total_amount_month = total_amount_month.reindex(index=month_order)
total_amount_month
```

```
[28]:
```

month	total_amount
January	31735.25
February	28937.89
March	33085.89
April	32012.54
May	33828.58
June	32920.52
July	26617.64
August	27759.56

September	28206.38
October	33065.83
November	30800.44
December	31261.57

```
[29]: plt.figure(figsize=(12,7))
ax = sns.barplot(x=total_amount_month.index,
                y=total_amount_month['total_amount'])
plt.title('Total revenue by month', fontsize=16);
```



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

```
[30]: df['DOLocationID'].nunique()
```

```
[30]: 216
```

```
[31]: distance_by_dropoff = df.groupby('DOLocationID').mean()[['trip_distance']]

distance_by_dropoff = distance_by_dropoff.sort_values(by='trip_distance')
distance_by_dropoff
```

```
[31]:
```

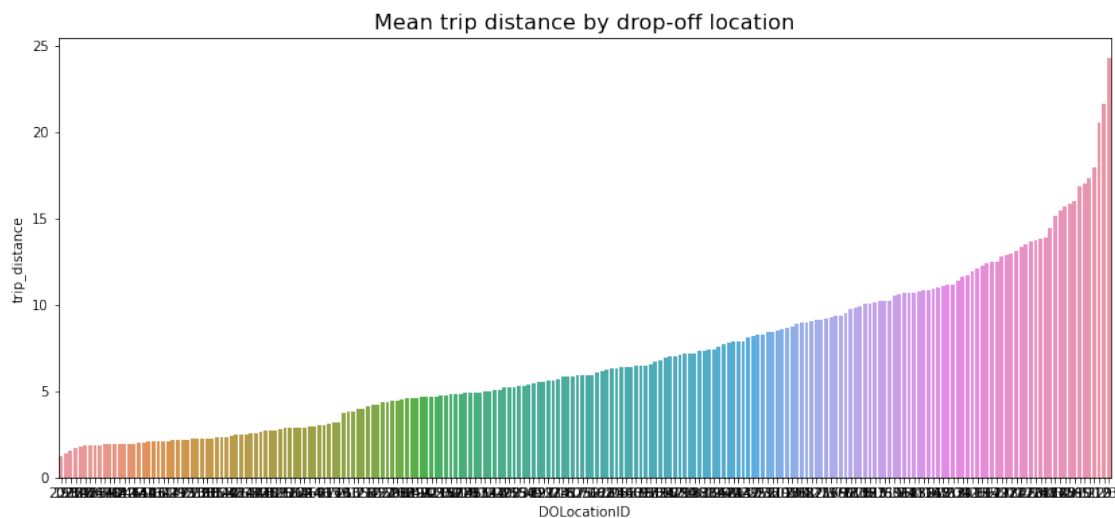
	trip_distance
DOLocationID	
207	1.200000
193	1.390556
237	1.555494
234	1.727806
137	1.818852
...	...
51	17.310000
11	17.945000
210	20.500000
29	21.650000
23	24.275000

[216 rows x 1 columns]

```
[32]: plt.figure(figsize=(14,6))
plt.title('Mean trip distance by drop-off location', fontsize=16);

sns.barplot(x=distance_by_dropoff.index,
            y=distance_by_dropoff['trip_distance'],
            order=distance_by_dropoff.index)
```

```
[32]: <matplotlib.axes._subplots.AxesSubplot at 0x7c90df7f5bd0>
```



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 for confirmation
```

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.

```
[ ]: #BONUS CONTENT for confirmation
```

To eliminate the spaces in the histogram 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.

```
[ ]: #BONUS CONTENT for confirmation
```

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 that there are 18 columns and not one has any missing value. Also looking at the barplots we can say that Monthly and daily ride counts generally aligns with monthly and daily revenues.

My other questions are why some trips have distances of zero. Also there are some outliers present, I want to know the specific reason behind them because that might provide a valuable insight

My client would likely want to know that the dataset includes trip start and end locations, allowing us to calculate the total distance traveled for each trip. This information could be valuable for optimizing routes or estimating fuel consumption, which may benefit the client's model..

```
[33]: df['trip_duration'] = (df['tpep_dropoff_datetime']-df['tpep_pickup_datetime'])
```

```
[ ]: df.head(10)
```

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 it involvest horough exploration of the dataset.It helps identify patterns, detect outliers, handle missing or inconsistent values, and ensure the data is properly prepared for modeling and deeper analysis.

Visualizations helped me understand key trends and potential anomalies in the data. They helped me identify several outliers, which will require careful consideration. ==> 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.