# Activity_ Course 5 Automatidata project lab

August 17, 2025

## 1 Automatidata project

**Course 5 - Regression Analysis: Simplify complex data relationships**

The data consulting firm Automatidata has recently hired you as the newest member of their data analytics team. Their newest client, the NYC Taxi and Limousine Commission (New York City TLC), wants the Automatidata team to build a multiple linear regression model to predict taxi fares using existing data that was collected over the course of a year. The team is getting closer to completing the project, having completed an initial plan of action, initial Python coding work, EDA, and A/B testing.

The Automatidata team has reviewed the results of the A/B testing. Now it's time to work on predicting the taxi fare amounts. You've impressed your Automatidata colleagues with your hard work and attention to detail. The data team believes that you are ready to build the regression model and update the client New York City TLC about your progress.

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

## 2 Course 5 End-of-course project: Build a multiple linear regression model

In this activity, you will build a multiple linear regression model. As you've learned, multiple linear regression helps you estimate the linear relationship between one continuous dependent variable and two or more independent variables. For data science professionals, this is a useful skill because it allows you to consider more than one variable against the variable you're measuring against. This opens the door for much more thorough and flexible analysis to be completed.

Completing this activity will help you practice planning out and buidling a multiple linear regression model 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 demostrate knowledge of EDA and a multiple linear regression model

**The goal** is to build a multiple linear regression model and evaluate the model *This activity has three parts:*

**Part 1:** EDA & Checking Model Assumptions * What are some purposes of EDA before constructing a multiple linear regression model?

**Part 2:** Model Building and evaluation * What resources do you find yourself using as you complete this stage?

**Part 3:** Interpreting Model Results

- What key insights emerged from your model(s)?
- What business recommendations do you propose based on the models built?

# 3   Build a multiple linear regression model

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.

### 3.0.1   Task 1. Imports and loading

Import the packages that you've learned are needed for building linear regression models.

```python
# Imports
# Packages for numerics + dataframes
import pandas as pd
import numpy as np

# Packages for visualization
import matplotlib.pyplot as plt
import seaborn as sns

# Packages for date conversions for calculating trip durations
from datetime import datetime
from datetime import date
from datetime import timedelta

# Packages for OLS, MLR, confusion matrix
from sklearn.preprocessing import StandardScaler
from sklearn.model_selection import train_test_split
import sklearn.metrics as metrics # For confusion matrix
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_absolute_error,r2_score,mean_squared_error
```

**Note:** `Pandas` is used to load the NYC TLC dataset. 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.

```
[3]: # Load dataset into dataframe
     df0=pd.read_csv("2017_Yellow_Taxi_Trip_Data.csv")
```

### 3.0.2 Task 2a. Explore data with EDA

Analyze and discover data, looking for correlations, missing data, outliers, and duplicates.

Start with `.shape` and `.info()`.

```
[4]: # Start with `.shape` and `.info()`
     ### YOUR CODE HERE ###
     df = df0.copy()
     df.shape
     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
```

Check for missing data and duplicates using `.isna()` and `.drop_duplicates()`.

```
[5]: # Check for missing data and duplicates using .isna() and .drop_duplicates()
     ### YOUR CODE HERE ###v
     df.describe()
```

[5]:
```
              Unnamed: 0        VendorID  passenger_count  trip_distance  \
count      2.269900e+04    22699.000000     22699.000000   22699.000000
mean       5.675849e+07        1.556236         1.642319       2.913313
std        3.274493e+07        0.496838         1.285231       3.653171
min        1.212700e+04        1.000000         0.000000       0.000000
25%        2.852056e+07        1.000000         1.000000       0.990000
50%        5.673150e+07        2.000000         1.000000       1.610000
75%        8.537452e+07        2.000000         2.000000       3.060000
max        1.134863e+08        2.000000         6.000000      33.960000

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

               extra       mta_tax    tip_amount  tolls_amount  \
count   22699.000000  22699.000000  22699.000000  22699.000000
mean        0.333275      0.497445      1.835781      0.312542
std         0.463097      0.039465      2.800626      1.399212
min        -1.000000     -0.500000      0.000000      0.000000
25%         0.000000      0.500000      0.000000      0.000000
50%         0.000000      0.500000      1.350000      0.000000
75%         0.500000      0.500000      2.450000      0.000000
max         4.500000      0.500000    200.000000     19.100000

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

Use .describe().

[6]:
```
# Use .describe()
### YOUR CODE HERE ###
```

### 3.0.3 Task 2b. Convert pickup & dropoff columns to datetime

```python
[7]:  #Convert datetime columns to datetime
      # Display data types of `tpep_pickup_datetime`, `tpep_dropoff_datetime`
      print('Data type of tpep_pickup_datetime:', df['tpep_pickup_datetime'].dtype)
      print('Data type of tpep_dropoff_datetime:', df['tpep_dropoff_datetime'].dtype)

      # Convert `tpep_pickup_datetime` to datetime format
      df['tpep_pickup_datetime'] = pd.to_datetime(df['tpep_pickup_datetime'],␣
       ↪format='%m/%d/%Y %I:%M:%S %p')

      # Convert `tpep_dropoff_datetime` to datetime format
      df['tpep_dropoff_datetime'] = pd.to_datetime(df['tpep_dropoff_datetime'],␣
       ↪format='%m/%d/%Y %I:%M:%S %p')

      # Display data types of `tpep_pickup_datetime`, `tpep_dropoff_datetime`
      print('Data type of tpep_pickup_datetime:', df['tpep_pickup_datetime'].dtype)
      print('Data type of tpep_dropoff_datetime:', df['tpep_dropoff_datetime'].dtype)

      df.head(3)
```

```
Data type of tpep_pickup_datetime: object
Data type of tpep_dropoff_datetime: object
Data type of tpep_pickup_datetime: datetime64[ns]
Data type of tpep_dropoff_datetime: datetime64[ns]
```

```
[7]:    Unnamed: 0  VendorID tpep_pickup_datetime tpep_dropoff_datetime  \
     0    24870114         2  2017-03-25 08:55:43   2017-03-25 09:09:47
     1    35634249         1  2017-04-11 14:53:28   2017-04-11 15:19:58
     2   106203690         1  2017-12-15 07:26:56   2017-12-15 07:34:08

        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

        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

        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
```

### 3.0.4  Task 2c.  Create duration column

Create a new column called `duration` that represents the total number of minutes that each taxi ride took.

```
[8]: df['duration'] = (df['tpep_dropoff_datetime'] - df['tpep_pickup_datetime'])/np.
     ↪timedelta64(1,'m')
```

### 3.0.5  Outliers

Call `df.info()` to inspect the columns and decide which ones to check for outliers.
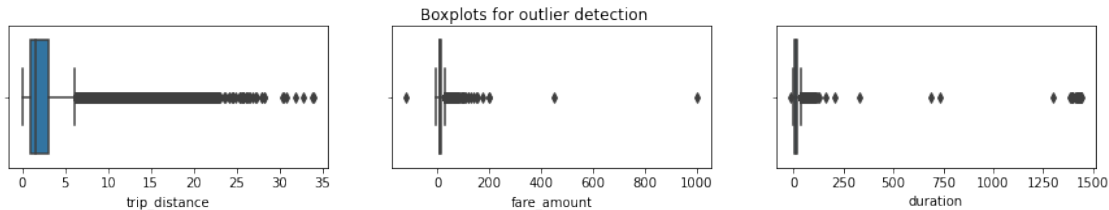
```
[9]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 22699 entries, 0 to 22698
Data columns (total 19 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  datetime64[ns]
 3   tpep_dropoff_datetime  22699 non-null  datetime64[ns]
 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
 18  duration               22699 non-null  float64
dtypes: datetime64[ns](2), float64(9), int64(7), object(1)
memory usage: 3.3+ MB
```

Keeping in mind that many of the features will not be used to fit your model, the most important columns to check for outliers are likely to be: * `trip_distance` * `fare_amount` * `duration`

### 3.0.6  Task 2d.  Box plots

Plot a box plot for each feature: `trip_distance`, `fare_amount`, `duration`.

```
[10]:  fig, axes = plt.subplots(1, 3, figsize=(15, 2))
       fig.suptitle('Boxplots for outlier detection')
       sns.boxplot(ax=axes[0], x=df['trip_distance'])
       sns.boxplot(ax=axes[1], x=df['fare_amount'])
       sns.boxplot(ax=axes[2], x=df['duration'])
       plt.show();
```



Boxplots for outlier detection

### 3.0.7 Task 2e. Imputations

**trip_distance outliers**   You know from the summary statistics that there are trip distances of 0. Are these reflective of erroneous data, or are they very short trips that get rounded down?

To check, sort the column values, eliminate duplicates, and inspect the least 10 values. Are they rounded values or precise values?

```
[11]:  sorted(set(df['trip_distance']))[:10]
```

```
[11]:  [0.0, 0.01, 0.02, 0.03, 0.04, 0.05, 0.06, 0.07, 0.08, 0.09]
```

The distances are captured with a high degree of precision. However, it might be possible for trips to have distances of zero if a passenger summoned a taxi and then changed their mind. Besides, are there enough zero values in the data to pose a problem?

Calculate the count of rides where the `trip_distance` is zero.

**fare_amount outliers   Question:** What do you notice about the values in the `fare_amount` column?

Impute values less than $0 with 0.

```
[12]:  df.loc[df['fare_amount'] < 0, 'fare_amount'] = 0
       df['fare_amount'].min()
```

```
[12]:  0.0
```

Now impute the maximum value as `Q3 + (6 * IQR)`.

```
[13]:  def outlier_imputer(column_list, iqr_factor):
           '''
```

```
    Impute upper-limit values in specified columns based on their interquartile␣
→range.

    Arguments:
        column_list: A list of columns to iterate over
        iqr_factor: A number representing x in the formula:
                    Q3 + (x * IQR). Used to determine maximum threshold,
                    beyond which a point is considered an outlier.

    The IQR is computed for each column in column_list and values exceeding
    the upper threshold for each column are imputed with the upper threshold␣
→value.
    '''
    for col in column_list:
        # Reassign minimum to zero
        df.loc[df[col] < 0, col] = 0

        # Calculate upper threshold
        q1 = df[col].quantile(0.25)
        q3 = df[col].quantile(0.75)
        iqr = q3 - q1
        upper_threshold = q3 + (iqr_factor * iqr)
        print(col)
        print('q3:', q3)
        print('upper_threshold:', upper_threshold)

        # Reassign values > threshold to threshold
        df.loc[df[col] > upper_threshold, col] = upper_threshold
        print(df[col].describe())
        print()
```

**duration outliers**

```
[14]: outlier_imputer(['fare_amount'], 6)
```

```
fare_amount
q3: 14.5
upper_threshold: 62.5
count    22699.000000
mean        12.897913
std         10.541137
min          0.000000
25%          6.500000
50%          9.500000
75%         14.500000
max         62.500000
```

```
Name: fare_amount, dtype: float64
```

The `duration` column has problematic values at both the lower and upper extremities.

- **Low values:** There should be no values that represent negative time. Impute all negative durations with 0.

- **High values:** Impute high values the same way you imputed the high-end outliers for fares: `Q3 + (6 * IQR)`.

```
[15]:  # Impute a 0 for any negative values
       ### YOUR CODE HERE ###
       df['duration'].describe()
```

```
[15]: count    22699.000000
      mean        17.013777
      std         61.996482
      min        -16.983333
      25%          6.650000
      50%         11.183333
      75%         18.383333
      max       1439.550000
      Name: duration, dtype: float64
```

```
[16]:  # Impute the high outliers
       ### YOUR CODE HERE ###df.loc[df['duration'] < 0, 'duration'] = 0
       df['duration'].min()
```

```
[16]: -16.983333333333334
```

```
[17]:  outlier_imputer(['duration'], 6)
```

```
duration
q3: 18.383333333333333
upper_threshold: 88.78333333333333
count    22699.000000
mean        14.460555
std         11.947043
min          0.000000
25%          6.650000
50%         11.183333
75%         18.383333
max         88.783333
Name: duration, dtype: float64
```

### 3.0.8 Task 3a. Feature engineering

```
[18]: df['pickup_dropoff'] = df['PULocationID'].astype(str) + ' ' +␣
       ↪df['DOLocationID'].astype(str)
      df['pickup_dropoff'].head(2)
```

```
[18]: 0      100 231
      1       186 43
      Name: pickup_dropoff, dtype: object
```

Now, use a `groupby()` statement to group each row by the new `pickup_dropoff` column, compute the mean, and capture the values only in the `trip_distance` column. Assign the results to a variable named `grouped`.

```
[19]: grouped = df.groupby('pickup_dropoff').
       ↪mean(numeric_only=True)[['trip_distance']]
      grouped[:5]
```

```
[19]:                  trip_distance
      pickup_dropoff
      1 1                   2.433333
      10 148               15.700000
      100 1                16.890000
      100 100               0.253333
      100 107               1.180000
```

`grouped` is an object of the `DataFrame` class.

1. Convert it to a dictionary using the `to_dict()` method. Assign the results to a variable called `grouped_dict`. This will result in a dictionary with a key of `trip_distance` whose values are another dictionary. The inner dictionary's keys are pickup/dropoff points and its values are mean distances. This is the information you want.

Example:
`grouped_dict = {'trip_distance': {'A B': 1.25, 'C D': 2, 'D C': 3}`

2. Reassign the `grouped_dict` dictionary so it contains only the inner dictionary. In other words, get rid of `trip_distance` as a key, so:

Example:
`grouped_dict = {'A B': 1.25, 'C D': 2, 'D C': 3}`

```
[20]: grouped_dict = grouped.to_dict()

      # 2. Reassign to only contain the inner dictionary
      grouped_dict = grouped_dict['trip_distance']
```

```
[21]: df['mean_distance'] = df['pickup_dropoff']
```

```
# 2. Map `grouped_dict` to the `mean_distance` column
df['mean_distance'] = df['mean_distance'].map(grouped_dict)

# Confirm that it worked
df[(df['PULocationID']==100) & (df['DOLocationID']==231)][['mean_distance']]
```

[21]:
```
        mean_distance
0            3.521667
4909         3.521667
16636        3.521667
18134        3.521667
19761        3.521667
20581        3.521667
```

**Create `mean_duration` column**  Repeat the process used to create the `mean_distance` column to create a `mean_duration` column.

[22]:
```
grouped = df.groupby('pickup_dropoff').mean(numeric_only=True)[['duration']]
grouped

# Create a dictionary where keys are unique pickup_dropoffs and values are
# mean trip duration for all trips with those pickup_dropoff combos
grouped_dict = grouped.to_dict()
grouped_dict = grouped_dict['duration']

df['mean_duration'] = df['pickup_dropoff']
df['mean_duration'] = df['mean_duration'].map(grouped_dict)

# Confirm that it worked
df[(df['PULocationID']==100) & (df['DOLocationID']==231)][['mean_duration']]
```

[22]:
```
        mean_duration
0           22.847222
4909        22.847222
16636       22.847222
18134       22.847222
19761       22.847222
20581       22.847222
```

**Create `day` and `month` columns**  Create two new columns, `day` (name of day) and `month` (name of month) by extracting the relevant information from the `tpep_pickup_datetime` column.

[23]:
```
# Create 'day' col
df['day'] = df['tpep_pickup_datetime'].dt.day_name().str.lower()

# Create 'month' col
```

```
df['month'] = df['tpep_pickup_datetime'].dt.strftime('%b').str.lower()
```

**Create rush_hour column**   Define rush hour as: * Any weekday (not Saturday or Sunday) AND
* Either from 06:00–10:00 or from 16:00–20:00

Create a binary `rush_hour` column that contains a 1 if the ride was during rush hour and a 0 if it
was not.

```
[24]: # Create 'rush_hour' col
      df['rush_hour'] = df['tpep_pickup_datetime'].dt.hour

      # If day is Saturday or Sunday, impute 0 in `rush_hour` column
      df.loc[df['day'].isin(['saturday', 'sunday']), 'rush_hour'] = 0
```

```
[25]: def rush_hourizer(hour):
          if 6 <= hour['rush_hour'] < 10:
              val = 1
          elif 16 <= hour['rush_hour'] < 20:
              val = 1
          else:
              val = 0
          return val
```

```
[26]: df.loc[(df.day != 'saturday') & (df.day != 'sunday'), 'rush_hour'] = df.
      ↪apply(rush_hourizer, axis=1)
      df.head()
```

```
[26]:    Unnamed: 0  VendorID tpep_pickup_datetime tpep_dropoff_datetime  \
      0    24870114         2  2017-03-25 08:55:43   2017-03-25 09:09:47
      1    35634249         1  2017-04-11 14:53:28   2017-04-11 15:19:58
      2   106203690         1  2017-12-15 07:26:56   2017-12-15 07:34:08
      3    38942136         2  2017-05-07 13:17:59   2017-05-07 13:48:14
      4    30841670         2  2017-04-15 23:32:20   2017-04-15 23:49:03

         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

         PULocationID  DOLocationID  …  tolls_amount  improvement_surcharge  \
      0           100           231  …           0.0                    0.3
      1           186            43  …           0.0                    0.3
      2           262           236  …           0.0                    0.3
      3           188            97  …           0.0                    0.3
      4             4           112  …           0.0                    0.3
```

```
     total_amount   duration  pickup_dropoff  mean_distance  mean_duration  \
0           16.56  14.066667             100 231      3.521667      22.847222
1           20.80  26.500000             186 43       3.108889      24.470370
2            8.75   7.200000             262 236      0.881429       7.250000
3           27.69  30.250000             188 97       3.700000      30.250000
4           17.80  16.716667               4 112      4.435000      14.616667

        day  month rush_hour
0  saturday    mar         0
1   tuesday    apr         0
2    friday    dec         1
3    sunday    may         0
4  saturday    apr         0

[5 rows x 25 columns]
```
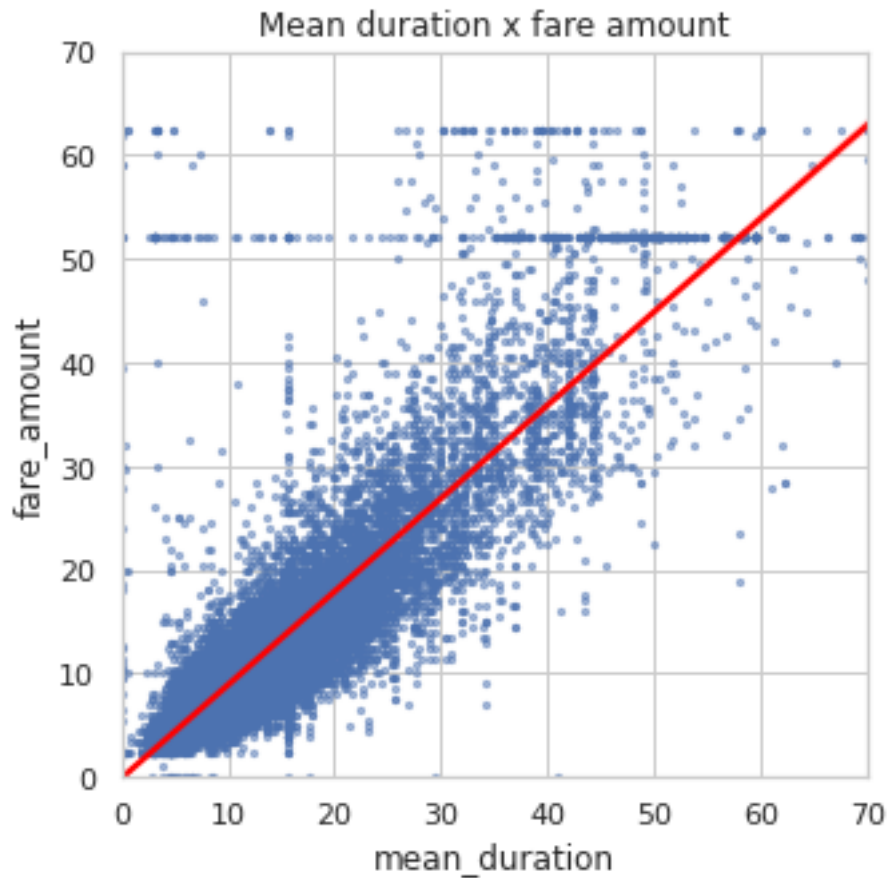
### 3.0.9  Task 4. Scatter plot

Create a scatterplot to visualize the relationship between `mean_duration` and `fare_amount`.

```python
[27]: sns.set(style='whitegrid')
      f = plt.figure()
      f.set_figwidth(5)
      f.set_figheight(5)
      sns.regplot(x=df['mean_duration'], y=df['fare_amount'],
                  scatter_kws={'alpha':0.5, 's':5},
                  line_kws={'color':'red'})
      plt.ylim(0, 70)
      plt.xlim(0, 70)
      plt.title('Mean duration x fare amount')
      plt.show()
```

Mean duration x fare amount

The `mean_duration` variable correlates with the target variable.

Check the value of the rides in the second horizontal line in the scatter plot.

```
[28]: df[df['fare_amount'] > 50]['fare_amount'].value_counts().head()
```

```
[28]: 52.0    514
      62.5     84
      59.0      9
      50.5      9
      57.5      8
      Name: fare_amount, dtype: int64
```

### 3.0.10 Task 5. Isolate modeling variables

Drop features that are redundant, irrelevant, or that will not be available in a deployed environment.

```
[29]: df2 = df.copy()
```

```
df2 = df2.drop(['Unnamed: 0', 'tpep_dropoff_datetime', 'tpep_pickup_datetime',
                'trip_distance', 'RatecodeID', 'store_and_fwd_flag',␣
↪'PULocationID', 'DOLocationID',
                'payment_type', 'extra', 'mta_tax', 'tip_amount',␣
↪'tolls_amount', 'improvement_surcharge',
                'total_amount', 'tpep_dropoff_datetime', 'tpep_pickup_datetime',␣
↪'duration',
                'pickup_dropoff', 'day', 'month'
                ], axis=1)

df2.info()
```
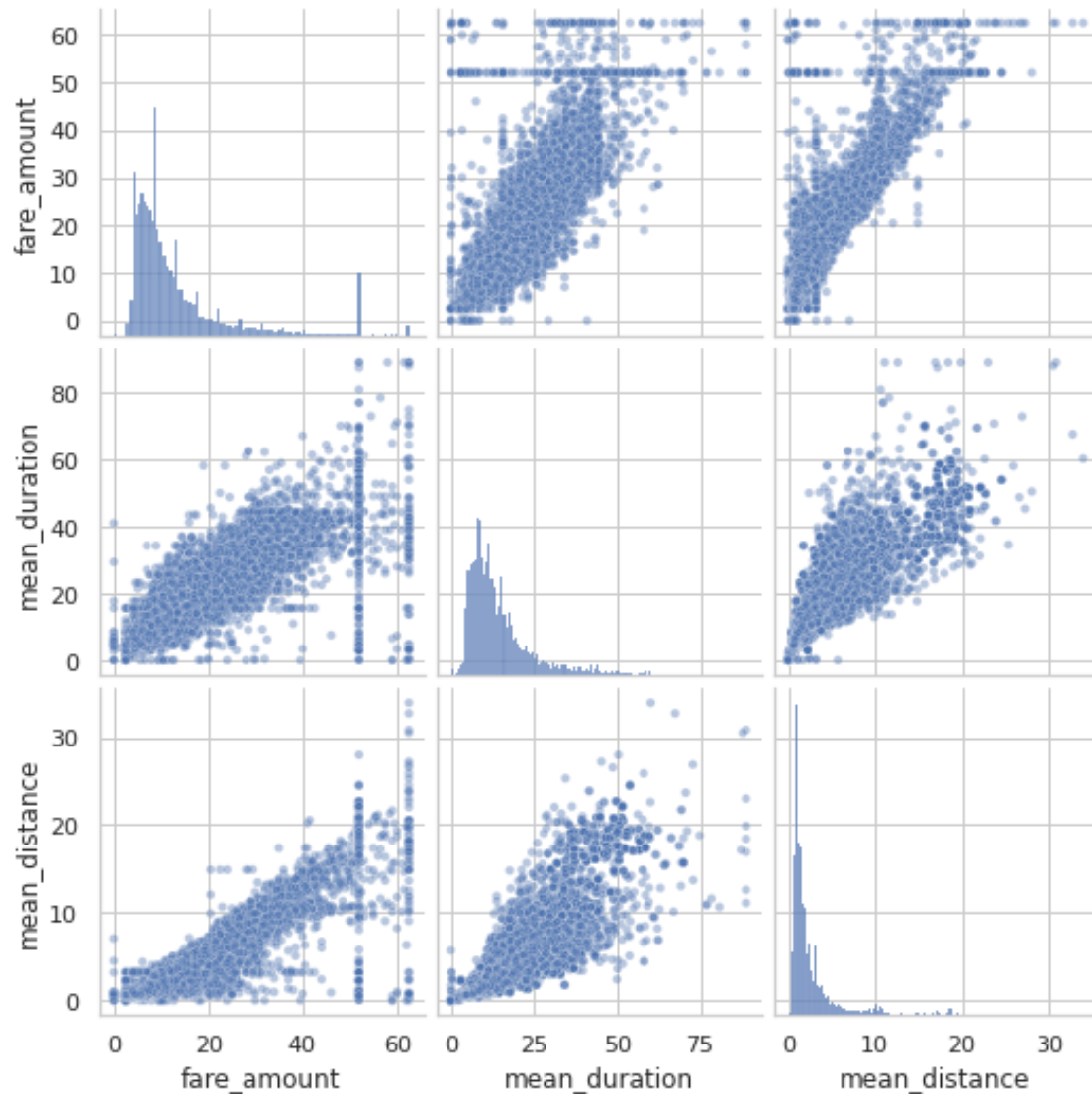
```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 22699 entries, 0 to 22698
Data columns (total 6 columns):
 #   Column           Non-Null Count  Dtype
---  ------           --------------  -----
 0   VendorID         22699 non-null  int64
 1   passenger_count  22699 non-null  int64
 2   fare_amount      22699 non-null  float64
 3   mean_distance    22699 non-null  float64
 4   mean_duration    22699 non-null  float64
 5   rush_hour        22699 non-null  int64
dtypes: float64(3), int64(3)
memory usage: 1.0 MB
```

### 3.0.11 Task 6. Pair plot

Create a pairplot to visualize pairwise relationships between `fare_amount`, `mean_duration`, and `mean_distance`.

```
[30]: sns.pairplot(df2[['fare_amount', 'mean_duration', 'mean_distance']],
                   plot_kws={'alpha':0.4, 'size':5},
                   );
```
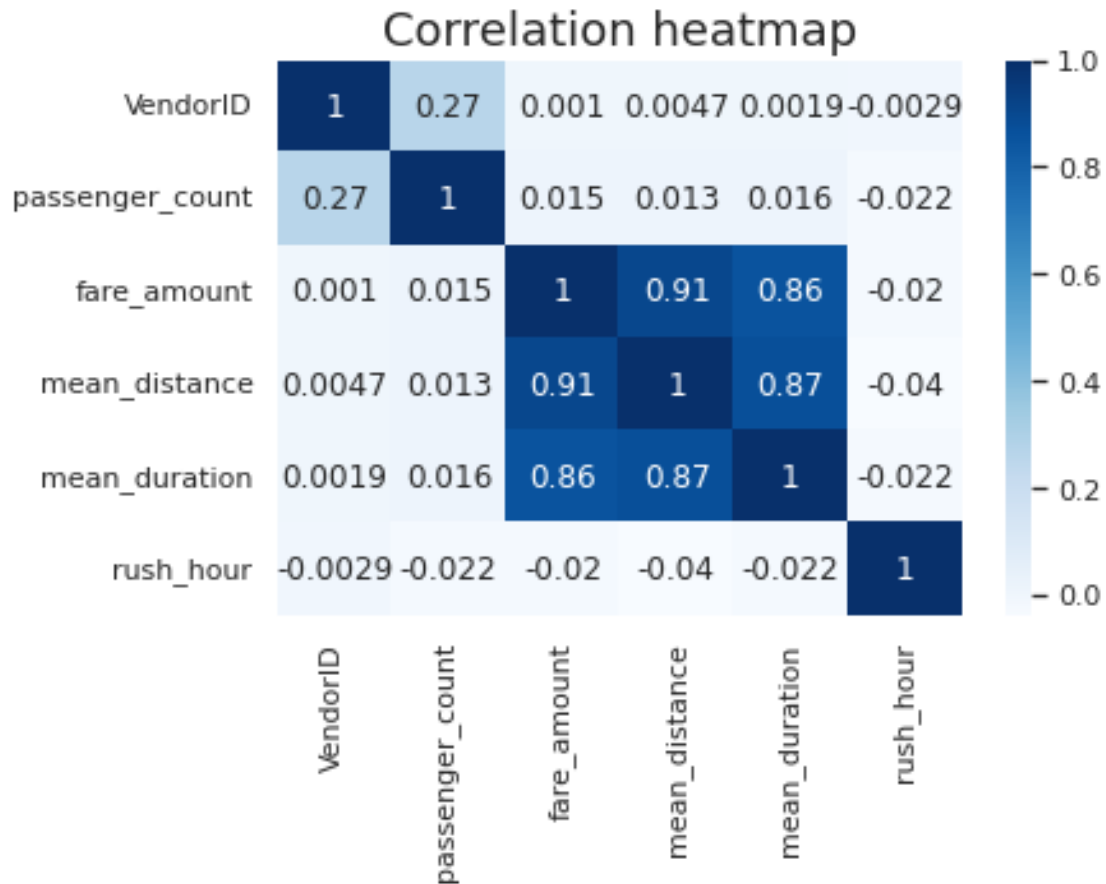
These variables all show linear correlation with each other. Investigate this further.

### 3.0.12   Task 7. Identify correlations

Visualize a correlation heatmap of the data.

```
[31]:  sns.heatmap(df2.corr(method='pearson'), annot=True, cmap='Blues')
       plt.title('Correlation heatmap',
                 fontsize=18)
       plt.show()
```

## Correlation heatmap

| | VendorID | passenger_count | fare_amount | mean_distance | mean_duration | rush_hour |
|---|---|---|---|---|---|---|
| VendorID | 1 | 0.27 | 0.001 | 0.0047 | 0.0019 | -0.0029 |
| passenger_count | 0.27 | 1 | 0.015 | 0.013 | 0.016 | -0.022 |
| fare_amount | 0.001 | 0.015 | 1 | 0.91 | 0.86 | -0.02 |
| mean_distance | 0.0047 | 0.013 | 0.91 | 1 | 0.87 | -0.04 |
| mean_duration | 0.0019 | 0.016 | 0.86 | 0.87 | 1 | -0.022 |
| rush_hour | -0.0029 | -0.022 | -0.02 | -0.04 | -0.022 | 1 |

### 3.0.13   Task 8a. Split data into outcome variable and features

Set your X and y variables. X represents the features and y represents the outcome (target) variable.

```
[32]:  X = df2.drop(columns=['fare_amount'])

       # Set y variable
       y = df2[['fare_amount']]

       # Display first few rows
       X.head()
```

```
[32]:     VendorID  passenger_count  mean_distance  mean_duration  rush_hour
       0         2                6       3.521667      22.847222          0
       1         1                1       3.108889      24.470370          0
       2         1                1       0.881429       7.250000          1
       3         2                1       3.700000      30.250000          0
       4         2                1       4.435000      14.616667          0
```

### 3.0.14 Task 8b. Pre-process data

Dummy encode categorical variables

```
[33]: X['VendorID'] = X['VendorID'].astype(str)

      # Get dummies
      X = pd.get_dummies(X, drop_first=True)
      X.head()
```

```
[33]:    passenger_count  mean_distance  mean_duration  rush_hour  VendorID_2
      0                6       3.521667      22.847222          0           1
      1                1       3.108889      24.470370          0           0
      2                1       0.881429       7.250000          1           0
      3                1       3.700000      30.250000          0           1
      4                1       4.435000      14.616667          0           1
```

### 3.0.15 Split data into training and test sets

Create training and testing sets. The test set should contain 20% of the total samples. Set
random_state=0.

```
[34]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,␣
      ↪random_state=0)
```

### 3.0.16 Standardize the data

Use `StandardScaler()`, `fit()`, and `transform()` to standardize the `X_train` variables. Assign
the results to a variable called `X_train_scaled`.

```
[35]: scaler = StandardScaler().fit(X_train)
      X_train_scaled = scaler.transform(X_train)
```

### 3.0.17 Fit the model

Instantiate your model and fit it to the training data.

```
[36]: lr=LinearRegression()
      lr.fit(X_train_scaled, y_train)
```

```
[36]: LinearRegression(copy_X=True, fit_intercept=True, n_jobs=None, normalize=False)
```

### 3.0.18   Task 8c.  Evaluate model

### 3.0.19   Train data

Evaluate your model performance by calculating the residual sum of squares and the explained variance score (R^2).  Calculate the Mean Absolute Error, Mean Squared Error, and the Root Mean Squared Error.

```
[37]: r_sq = lr.score(X_train_scaled, y_train)
      print('Coefficient of determination:', r_sq)
      y_pred_train = lr.predict(X_train_scaled)
      print('R^2:', r2_score(y_train, y_pred_train))
      print('MAE:', mean_absolute_error(y_train, y_pred_train))
      print('MSE:', mean_squared_error(y_train, y_pred_train))
      print('RMSE:',np.sqrt(mean_squared_error(y_train, y_pred_train)))
```

```
Coefficient of determination: 0.8398434585044773
R^2: 0.8398434585044773
MAE: 2.186666416775414
MSE: 17.88973296349268
RMSE: 4.229625629236313
```

### 3.0.20   Test data

Calculate the same metrics on the test data. Remember to scale the `X_test` data using the scaler that was fit to the training data. Do not refit the scaler to the testing data, just transform it. Call the results `X_test_scaled`.

```
[38]: X_test_scaled = scaler.transform(X_test)
```

```
[39]: r_sq_test = lr.score(X_test_scaled, y_test)
      print('Coefficient of determination:', r_sq_test)
      y_pred_test = lr.predict(X_test_scaled)
      print('R^2:', r2_score(y_test, y_pred_test))
      print('MAE:', mean_absolute_error(y_test,y_pred_test))
      print('MSE:', mean_squared_error(y_test, y_pred_test))
      print('RMSE:',np.sqrt(mean_squared_error(y_test, y_pred_test)))
```

```
Coefficient of determination: 0.8682583641795454
R^2: 0.8682583641795454
MAE: 2.1336549840593864
MSE: 14.326454156998944
RMSE: 3.785030271609323
```

### 3.0.21 Task 9a. Results

Use the code cell below to get `actual`,`predicted`, and `residual` for the testing set, and store them as columns in a `results` dataframe.

```
[40]: results = pd.DataFrame(data={'actual': y_test['fare_amount'],
                                   'predicted': y_pred_test.ravel()})
      results['residual'] = results['actual'] - results['predicted']
      results.head()
```

```
[40]:        actual  predicted   residual
       5818    14.0  12.356503   1.643497
       18134   28.0  16.314595  11.685405
       4655     5.5   6.726789  -1.226789
       7378    15.5  16.227206  -0.727206
       13914    9.5  10.536408  -1.036408
```

```
[41]: coefficients = pd.DataFrame(lr.coef_, columns=X.columns)
      coefficients
```

```
[41]:    passenger_count  mean_distance  mean_duration  rush_hour  VendorID_2
       0         0.030825       7.133867       2.812115   0.110233   -0.054373
```

```
[42]: print(X_train['mean_distance'].std())

      # 2. Divide the model coefficient by the standard deviation
      print(7.133867 / X_train['mean_distance'].std())
```

```
3.574812975256415
1.9955916713344426
```

Now you can make a more intuitive interpretation: for every 3.57 miles traveled, the fare increased by a mean of \$7.13. Or, reduced: for every 1 mile traveled, the fare increased by a mean of \$2.00.