

# Background on the Automatidata scenario

Automatidata works with its clients to transform their unused and stored data into useful solutions, such as performance dashboards, customer-facing tools, strategic business insights, and more. They specialize in identifying a client's business needs and utilizing their data to meet those business needs.

Automatidata is consulting for the New York City Taxi and Limousine Commission (TLC). New York City TLC is an agency responsible for licensing and regulating New York City's taxi cabs and for-hire vehicles. The agency has partnered with Automatidata to develop a regression model that predicts taxi and limousine ride durations based on location and time of day data that TLC has gathered.

The TLC data comes from over 200,000 taxi and limousine licensees, making approximately one million combined trips per day.

## Project background

Automatidata is near the end of the TLC project. The following tasks are needed at this stage of the project:

1. Determine the correct modeling approach
2. Build a regression model
3. Finish checking model assumptions
4. Evaluate the model
5. Interpret model results and summarize findings for stakeholders within TLC

## Automatidata Teammates

Udo Bankole, Director of Data Analysis

Deshawn Washington, Data Analysis Manager

Luana Rodriguez, Senior Data Analyst

Uli King, Senior Project Manager

## Stakeholders

Juliana Soto, Finance and Administration Department Head

Titus Nelson, Operations Manager

## Project goal

1. Build multiple linear regression model
2. Evaluate the model

```

In [1]: 1 # Imports
        2
        3 # Packages for numerics + dataframes
        4 import pandas as pd
        5 import numpy as np
        6
        7 # Packages for visualization
        8 import matplotlib.pyplot as plt
        9 import seaborn as sns
       10
       11 # Packages for date conversions for calculating trip durations
       12 from datetime import datetime
       13 from datetime import date
       14 from datetime import timedelta
       15
       16 # Packages for OLS, MLR, confusion matrix
       17 from sklearn.preprocessing import StandardScaler
       18 from sklearn.model_selection import train_test_split
       19 import sklearn.metrics as metrics # For confusion matrix
       20 from sklearn.linear_model import LinearRegression
       21 from sklearn.metrics import mean_absolute_error, r2_score, mean_squared_error
       22 lr=LinearRegression()

```

## Analyze pphase/EDA

```

In [3]: 1
        2 # import data
        3 df0=pd.read_csv("archive (3).zip")

```

```

In [5]: 1 df0.head(10) # first ten rows

```

```

Out[5]:

```

	Unnamed: 0	VendorID	tpep_pickup_datetime	tpep_dropoff_datetime	passenger_count	trip_distance	Ratecor
0	24870114	2	3/25/2017 8:55	3/25/2017 9:09	6	3.34	
1	35634249	1	4/11/2017 14:53	4/11/2017 15:19	1	1.80	
2	106203690	1	12/15/2017 7:26	12/15/2017 7:34	1	1.00	
3	38942136	2	5/7/2017 13:17	5/7/2017 13:48	1	3.70	
4	30841670	2	4/15/2017 23:32	4/15/2017 23:49	1	4.37	
5	23345809	2	3/25/2017 20:34	3/25/2017 20:42	6	2.30	
6	37660487	2	5/3/2017 19:04	5/3/2017 20:03	1	12.83	
7	69059411	2	8/15/2017 17:41	8/15/2017 18:03	1	2.98	
8	8433159	2	2/4/2017 16:17	2/4/2017 16:29	1	1.20	
9	95294817	1	11/10/2017 15:20	11/10/2017 15:40	1	1.60	

```
In [6]: 1 df0.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
```

```
In [7]: 1
2 df = df0.copy()
3
4 # Display the dataset's shape
5 print(df.shape)
6
7 # Display basic info about the dataset
8 df.info()
```

```
(22699, 18)
<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
```

```

In [9]: 1 # Create `trip_duration`
        2
        3 # Display data types of `tpep_dropoff_datetime`, `tpep_pickup_datetime`
        4 print("Data type of tpep_dropoff_datetime:", df["tpep_dropoff_datetime"].dtype)
        5 print("Data type of tpep_pickup_datetime:", df["tpep_pickup_datetime"].dtype)
        6
        7 # Convert `tpep_dropoff_datetime` to datetime format
        8 df["drop_off_converted"] = pd.to_datetime(df["tpep_dropoff_datetime"])
        9
       10 # Convert `tpep_pickup_datetime` to datetime format
       11 df["pick_up_converted"] = pd.to_datetime(df["tpep_pickup_datetime"])
       12
       13 # Display data types of `drop_off_converted`, `pick_up_converted`
       14 print("Data type of drop_off_converted:", df["drop_off_converted"].dtype)
       15 print("Data type of pick_up_converted:", df["pick_up_converted"].dtype)
       16
       17 # Compute `trip_duration`
       18 df["trip_duration"] = (df["drop_off_converted"] - df["pick_up_converted"])/np.timedelta64(1, 'ns')
       19
       20 # Display first ten rows of dataframe after adding the new columns
       21 df.head(10)

```

Data type of tpep\_dropoff\_datetime: object  
 Data type of tpep\_pickup\_datetime: object  
 Data type of drop\_off\_converted: datetime64[ns]  
 Data type of pick\_up\_converted: datetime64[ns]

Out[9]:

	Unnamed: 0	VendorID	tpep_pickup_datetime	tpep_dropoff_datetime	passenger_count	trip_distance	Ratecode
0	24870114	2	3/25/2017 8:55	3/25/2017 9:09	6	3.34	
1	35634249	1	4/11/2017 14:53	4/11/2017 15:19	1	1.80	
2	106203690	1	12/15/2017 7:26	12/15/2017 7:34	1	1.00	
3	38942136	2	5/7/2017 13:17	5/7/2017 13:48	1	3.70	
4	30841670	2	4/15/2017 23:32	4/15/2017 23:49	1	4.37	
5	23345809	2	3/25/2017 20:34	3/25/2017 20:42	6	2.30	
6	37660487	2	5/3/2017 19:04	5/3/2017 20:03	1	12.83	
7	69059411	2	8/15/2017 17:41	8/15/2017 18:03	1	2.98	
8	8433159	2	2/4/2017 16:17	2/4/2017 16:29	1	1.20	
9	95294817	1	11/10/2017 15:20	11/10/2017 15:40	1	1.60	

10 rows × 21 columns



```
In [11]: 1 # Check for missing data and duplicates
2
3 # Check for duplicates
4 print("Shape of dataframe:", df.shape)
5 print("Shape of dataframe with duplicates dropped:", df.drop_duplicates().shape)
6
7 # Check for missing values in dataframe
8 print("Total count of missing values:", df.isna().sum().sum())
9
10 # Display missing values per column in dataframe
11 print("Missing values per column:")
12 df.isna().sum()
```

Shape of dataframe: (22699, 21)

Shape of dataframe with duplicates dropped: (22699, 21)

Total count of missing values: 0

Missing values per column:

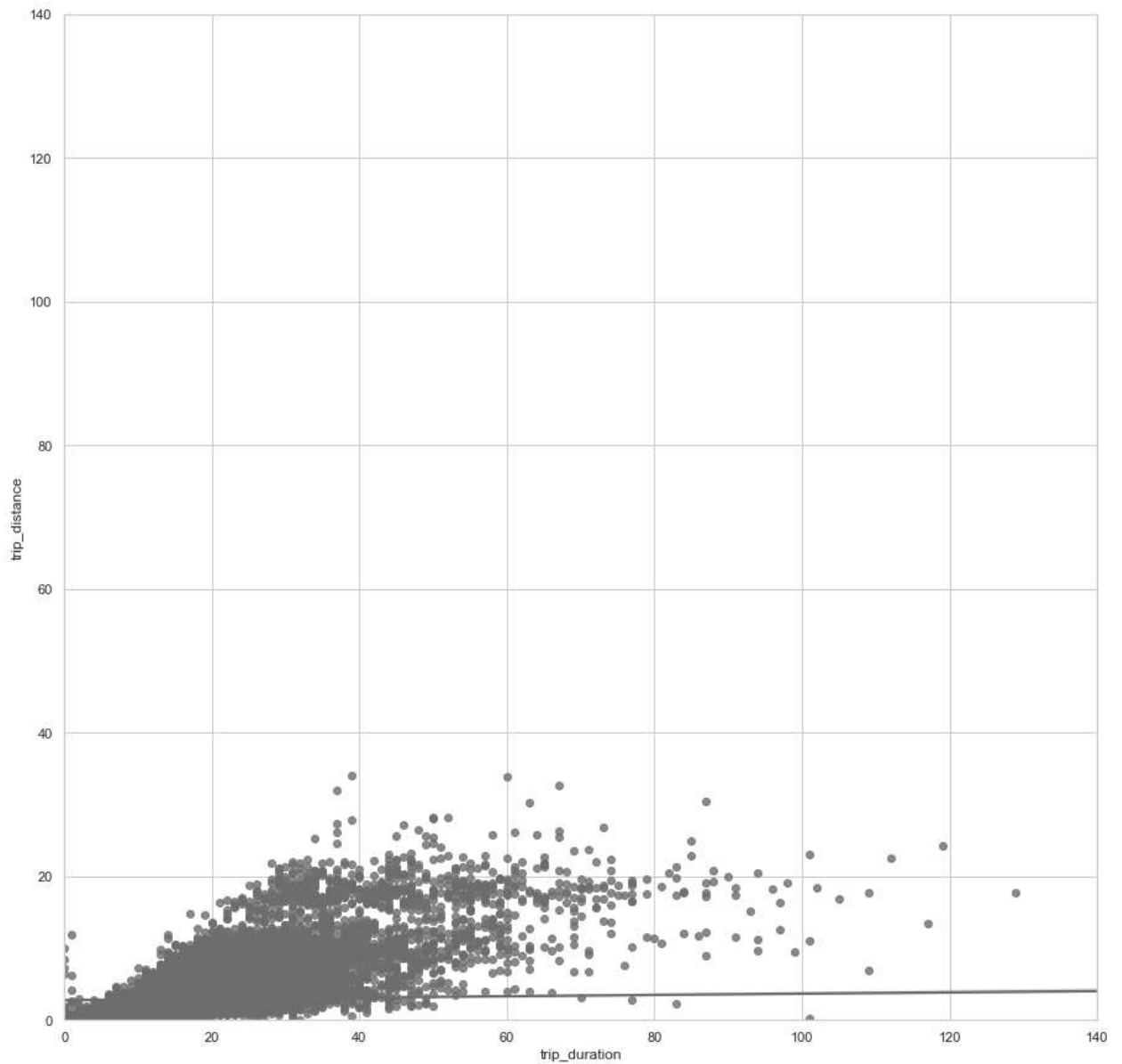
```
Out[11]: Unnamed: 0      0
VendorID      0
tpep_pickup_datetime      0
tpep_dropoff_datetime      0
passenger_count      0
trip_distance      0
RatecodeID      0
store_and_fwd_flag      0
PULocationID      0
DOLocationID      0
payment_type      0
fare_amount      0
extra      0
mta_tax      0
tip_amount      0
tolls_amount      0
improvement_surcharge      0
total_amount      0
drop_off_converted      0
pick_up_converted      0
trip_duration      0
dtype: int64
```

```
In [12]: 1 # Display descriptive stats about the data
2 df.describe()
```

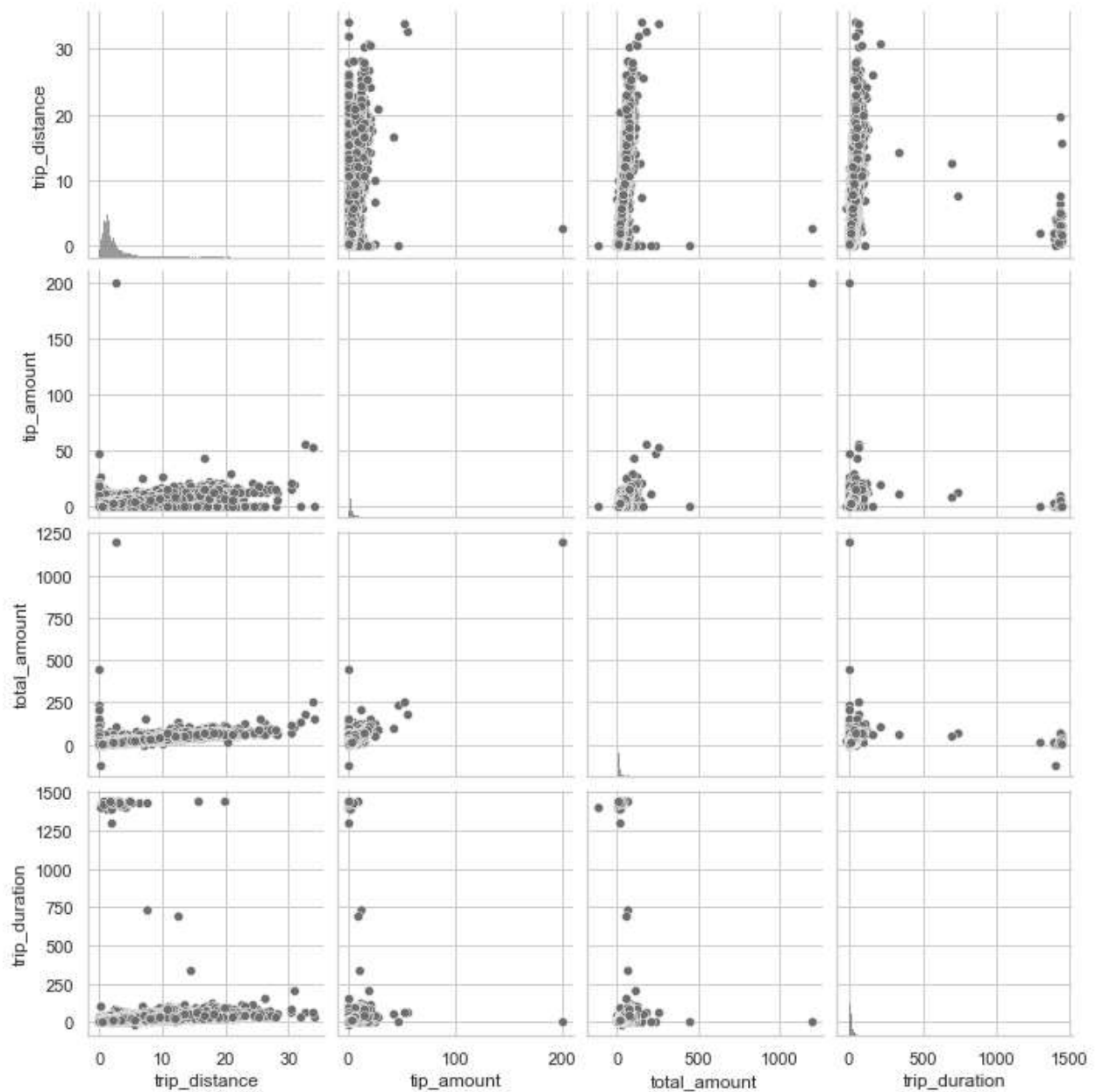
```
Out[12]:
```

	PULocationID	DOLocationID	payment_type	fare_amount	extra	mta_tax	tip_amount	tolls_amount	improvement_surcharge	total_amount	drop_off_converted	pick_up_converted	trip_duration
count	22699	22699	22699	22699	22699	22699	22699	22699	22699	22699	22699	22699	22699
mean	161.527997	70.139691	1.000000	13.243791	0.463097	0.039465	1.835781	0.312542	0.000000	19.100000	0.000000	0.000000	19.100000
std	123.530143	70.139691	0.496211	13.243791	0.463097	0.039465	1.835781	0.312542	0.000000	19.100000	0.000000	0.000000	19.100000
min	1.000000	1.000000	1.000000	-120.000000	-1.000000	-0.500000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
25%	161.527997	70.139691	1.000000	13.243791	0.463097	0.039465	1.835781	0.312542	0.000000	19.100000	0.000000	0.000000	19.100000
50%	161.527997	70.139691	1.000000	13.243791	0.463097	0.039465	1.835781	0.312542	0.000000	19.100000	0.000000	0.000000	19.100000
75%	161.527997	70.139691	1.000000	13.243791	0.463097	0.039465	1.835781	0.312542	0.000000	19.100000	0.000000	0.000000	19.100000
max	161.527997	70.139691	1.000000	13.243791	0.463097	0.039465	1.835781	0.312542	0.000000	19.100000	0.000000	0.000000	19.100000

```
In [13]: 1
2 # Create a scatter plot of trip_duration and trip_distance, with a line of best fit
3 sns.set(style='whitegrid')
4 f = plt.figure()
5 f.set_figwidth(15)
6 f.set_figheight(15)
7 sns.regplot(x=df["trip_duration"], y=df["trip_distance"])
8 plt.ylim(0, 140)
9 plt.xlim(0,140)
10 plt.show()
```



```
In [14]: 1 # Create a pairplot to visualize pairwise relationships between variables in the data
          2
          3 sns.pairplot(df[['trip_distance', 'tip_amount', 'total_amount', 'trip_duration']]);
```



**Address any outliers**

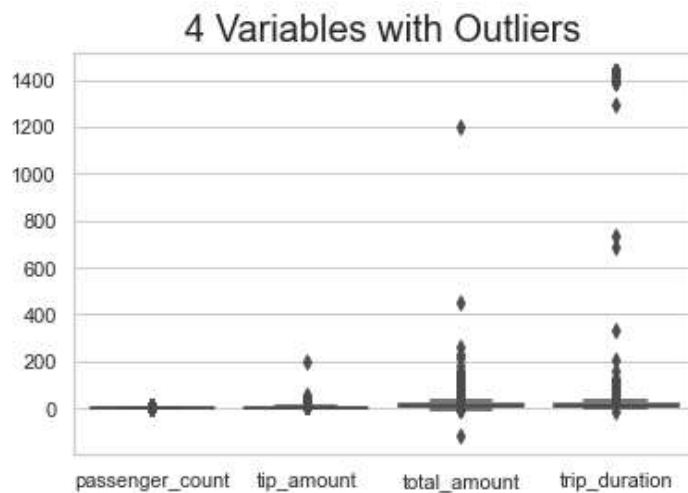


```

In [15]: 1 # Create boxplot to visualize the outliers
          2
          3 g = sns.boxplot(data=df[["passenger_count","tip_amount","total_amount", "trip_duration"]],
          4 g.set_title("4 Variables with Outliers",fontsize=20)

```

Out[15]: Text(0.5, 1.0, '4 Variables with Outliers')

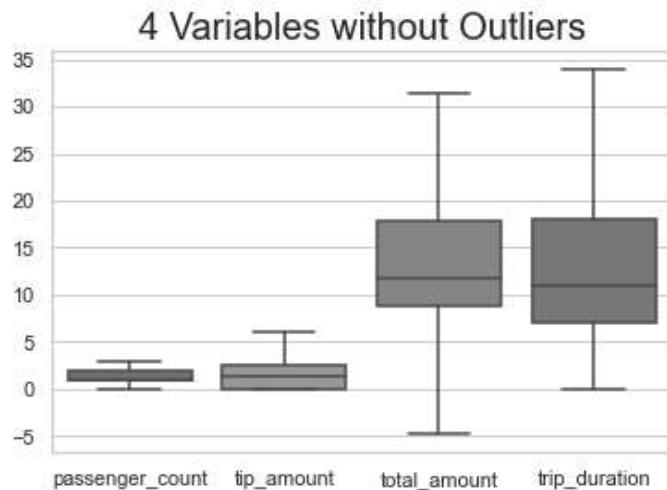


```

In [16]: 1 # Create boxplot to visualize distribution of data without outliers
          2
          3 g = sns.boxplot(data=df[["passenger_count","tip_amount","total_amount", "trip_duration"]],
          4 g.set_title("4 Variables without Outliers",fontsize=20)

```

Out[16]: Text(0.5, 1.0, '4 Variables without Outliers')



```

In [17]: 1 # Compute the 25th and 75th percentile values in `trip_duration`
          2 percentile25 = df["trip_duration"].quantile(0.25)
          3 percentile75 = df["trip_duration"].quantile(0.75)
          4
          5 # Compute the interquartile range for `trip_duration`
          6 iqr = percentile75 - percentile25
          7
          8 # Compute the upper limit for `trip_duration`
          9 upper_limit = percentile75 + 1.5 * iqr
          10 upper_limit

```

Out[17]: 34.5

```
In [18]: 1 # Remove outliers in `trip_duration`:
2 # Set values greater than the upper limit to the upper limit (approximately 36)
3 # Set values less than 0 to 0
4 df[df["trip_duration"] > 36] = 36
5 df[df["trip_duration"] < 0] = 0
6
7 df["trip_duration"].describe()
```

```
Out[18]: count    22699.000000
mean         13.663686
std           9.254188
min           0.000000
25%           7.000000
50%          11.000000
75%          18.000000
max          36.000000
Name: trip_duration, dtype: float64
```

```
In [19]: 1 # Compute the 25th and 75th percentile values in `total_amount`
2 percentile25 = df["total_amount"].quantile(0.25)
3 percentile75 = df["total_amount"].quantile(0.75)
4
5 # Compute the interquartile range for `total_amount`
6 iqr = percentile75 - percentile25
7
8 # Compute the upper limit for `total_amount`
9 upper_limit = percentile75 + 1.5 * iqr
10 upper_limit
```

```
Out[19]: 31.375
```

```
In [20]: 1 # Remove outliers in `total_amount`:
2 # Set values greater than the upper limit to the upper limit (approximately 32)
3 # Set values less than 0 to 0
4 df[df["total_amount"] > 32] = 32
5 df[df["total_amount"] < 0] = 0
6
7 # Display descriptive stats after removing outliers in `total_amount`
8 df["total_amount"].describe()
```

```
Out[20]: count    22699.000000
mean         14.478729
std           7.981224
min           0.000000
25%           8.750000
50%          11.800000
75%          17.800000
max          32.000000
Name: total_amount, dtype: float64
```

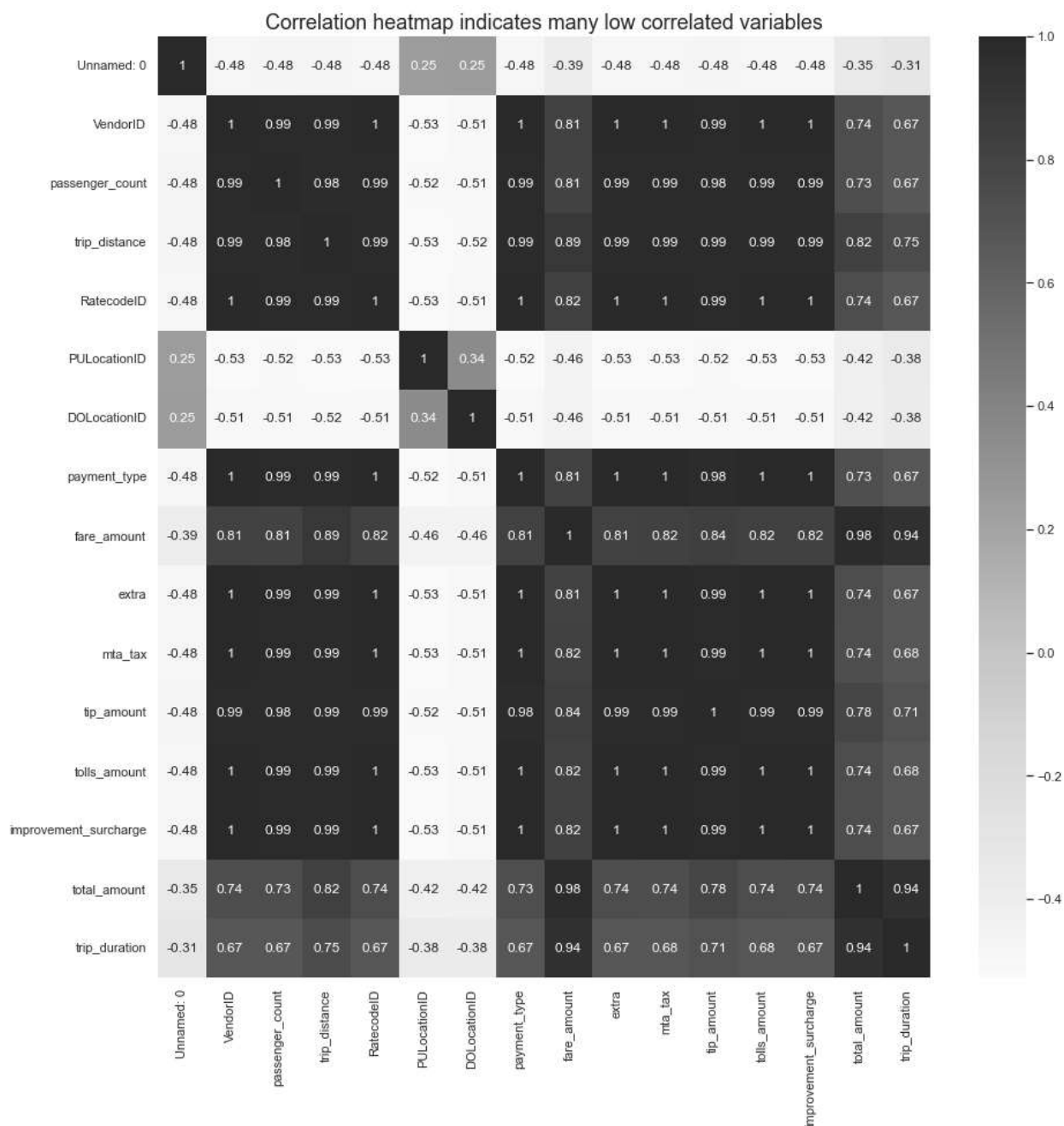
## Identify correlations

```
In [21]: 1
2 # Create correlation matrix containing pairwise correlation of columns, using pearson
3 df.corr(method="pearson")
```

Out[21]:

	Unnamed: 0	VendorID	passenger_count	trip_distance	RatecodeID	PULocationID	DOLocationID	payment_type	fare_amount	extra	mta_tax	tip_amount	tolls_amount	improvement_surcharge	total_amount	trip_duration
Unnamed: 0	1.000000	-0.480352	-0.476459	-0.476197	-0.481048	0.245243										
VendorID	-0.480352	1.000000	0.991809	0.986387	0.998622	-0.525293										
passenger_count	-0.476459	0.991809	1.000000	0.979355	0.991223	-0.521025										
trip_distance	-0.476197	0.986387	0.979355	1.000000	0.987438	-0.527592										
RatecodeID	-0.481048	0.998622	0.991223	0.987438	1.000000	-0.525036										
PULocationID	0.245243	-0.525293	-0.521025	-0.527592	-0.525036	1.000000										
DOLocationID	0.250318	-0.512225	-0.507229	-0.518649	-0.512051	0.343807	1.000000									
payment_type	-0.480212	0.997359	0.990115	0.986073	0.998661	-0.524588		1.000000								
fare_amount	-0.386649	0.814052	0.808990	0.887596	0.815310	-0.455797			1.000000							
extra	-0.481011	0.998006	0.990576	0.987081	0.999235	-0.525037				1.000000						
mta_tax	-0.481147	0.998684	0.991297	0.987574	0.999911	-0.525102					1.000000					
tip_amount	-0.476850	0.988978	0.981527	0.985181	0.990265	-0.520802						1.000000				
tolls_amount	-0.481163	0.998259	0.990888	0.987646	0.999529	-0.525127							1.000000			
improvement_surcharge	-0.481169	0.998683	0.991297	0.987569	0.999923	-0.525134								1.000000		
total_amount	-0.349756	0.736321	0.731639	0.821739	0.737534	-0.416848									1.000000	
trip_duration	-0.309033	0.673751	0.670005	0.747596	0.674577	-0.377790										1.000000

```
In [22]: 1 # Create correlation heatmap
2
3 plt.figure(figsize=(15,15))
4 sns.heatmap(df.corr(method="pearson"), annot=True, cmap="Blues") #cmap="crest")
5 plt.title("Correlation heatmap indicates many low correlated variables",
6           fontsize=18)
7 plt.show()
```



## Construct phase

```
In [23]: 1 # Set your Y and X variables
2
3 # Set Y variable
4 Y = df[["trip_duration"]]
5
6 # Remove the target column from the features
7 X = df.drop(columns="trip_duration")
8
9 # Display first few rows
10 X.head()
```

```
Out[23]:
```

	Unnamed: 0	VendorID	tpep_pickup_datetime	tpep_dropoff_datetime	passenger_count	trip_distance	RatecodeID
0	24870114	2	3/25/2017 8:55	3/25/2017 9:09	6	3.34	
1	35634249	1	4/11/2017 14:53	4/11/2017 15:19	1	1.80	
2	106203690	1	12/15/2017 7:26	12/15/2017 7:34	1	1.00	
3	38942136	2	5/7/2017 13:17	5/7/2017 13:48	1	3.70	
4	30841670	2	4/15/2017 23:32	4/15/2017 23:49	1	4.37	

```
In [ ]:
```

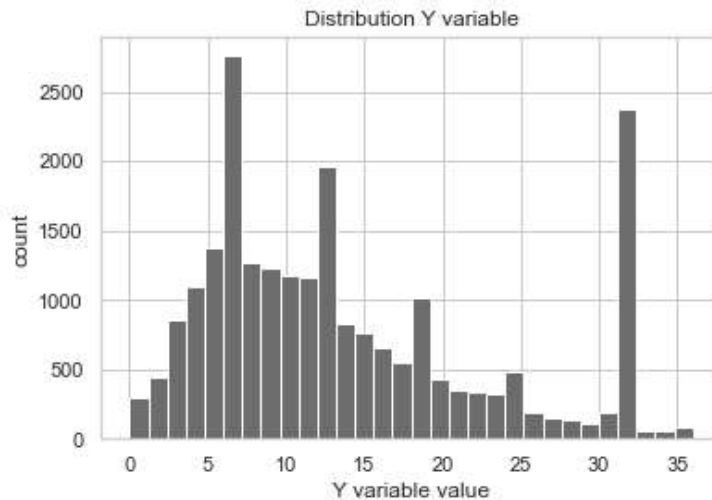
```
1
```

```
In [24]: 1 # columns to drop
2
3 columns_to_drop = ['tpep_pickup_datetime', 'tpep_dropoff_datetime',
4                   'store_and_fwd_flag', 'passenger_count', 'VendorID',
5                   'fare_amount', 'PULocationID', 'DOLocationID', 'total_amount',
6                   'drop_off_converted', 'pick_up_converted']
7 X = X.drop(columns_to_drop, axis=1)
8 X = X.loc[:, ~X.columns.str.contains("Unnamed")]
9 X.head()
```

```
Out[24]:
```

	trip_distance	RatecodeID	payment_type	extra	mta_tax	tip_amount	tolls_amount	improvement_surcharge
0	3.34	1	1	0.0	0.5	2.76	0.0	0.3
1	1.80	1	1	0.0	0.5	4.00	0.0	0.3
2	1.00	1	1	0.0	0.5	1.45	0.0	0.3
3	3.70	1	1	0.0	0.5	6.39	0.0	0.3
4	4.37	1	2	0.5	0.5	0.00	0.0	0.3

```
In [30]: 1 plt.hist(Y, bins=30)
2 plt.title("Distribution Y variable")
3 plt.xlabel("Y variable value")
4 plt.ylabel("count")
5 plt.show()
```



```
In [31]: 1 # Standardize the X variables
2 scaler = StandardScaler()
3 X_scaled = scaler.fit_transform(X)
4 print("X scaled:", X_scaled)
```

```
X scaled: [[-0.18035575 -0.33618447 -0.37233967 ... -0.18408591 -0.33725472
-0.33583702]
[-0.34804254 -0.33618447 -0.37233967 ... -0.05123384 -0.33725472
-0.33583702]
[-0.43515255 -0.33618447 -0.37233967 ... -0.32443769 -0.33725472
-0.33583702]
...
[-0.49830731 -0.33618447 -0.26440748 ... -0.4797889 -0.33725472
-0.33583702]
[-0.28706552 -0.33618447 -0.37233967 ... -0.297653 -0.33725472
-0.33583702]
[-0.31537628 -0.33618447 -0.37233967 ... -0.2280128 -0.33725472
-0.33583702]]
```

## Build model

```
In [32]: 1 # Create training and testing sets
2
3 X_train, X_test, Y_train, Y_test = train_test_split(X_scaled, Y, test_size=0.2, random_state=42)
```

```
In [33]: 1 # Build and fit your model to the training data
2 lr=LinearRegression()
3 lr.fit(X_train,Y_train)
```

```
Out[33]: LinearRegression()
```

## Evaluate model

```
In [34]: 1 # Evaluate the model performance on the training data
2 r_sq = lr.score(X_train, Y_train)
3 print("Coefficient of determination:", r_sq)
4 Y_pred = lr.predict(X_train)
5 print("R^2:", r2_score(Y_train, Y_pred))
6 print("MAE:", mean_absolute_error(Y_train, Y_pred))
7 print("RMSE:", np.sqrt(mean_squared_error(Y_train, Y_pred)))
```

Coefficient of determination: 0.7429168111442532

R^2: 0.7429168111442532

MAE: 3.2405134739474537

RMSE: 4.61484392596783

```
In [35]: 1 # Evaluate the model performance on the testing data
2 r_sq_test = lr.score(X_test, Y_test)
3 print("Coefficient of determination:", r_sq_test)
4 Y_pred_test = lr.predict(X_test)
5 print("R^2:", r2_score(Y_test, Y_pred_test))
6 print("MAE:", mean_absolute_error(Y_test, Y_pred_test))
7 print("RMSE:", np.sqrt(mean_squared_error(Y_test, Y_pred_test)))
```

Coefficient of determination: 0.736833952588452

R^2: 0.736833952588452

MAE: 3.253815127917235

RMSE: 4.632230695076403

## Execute phase

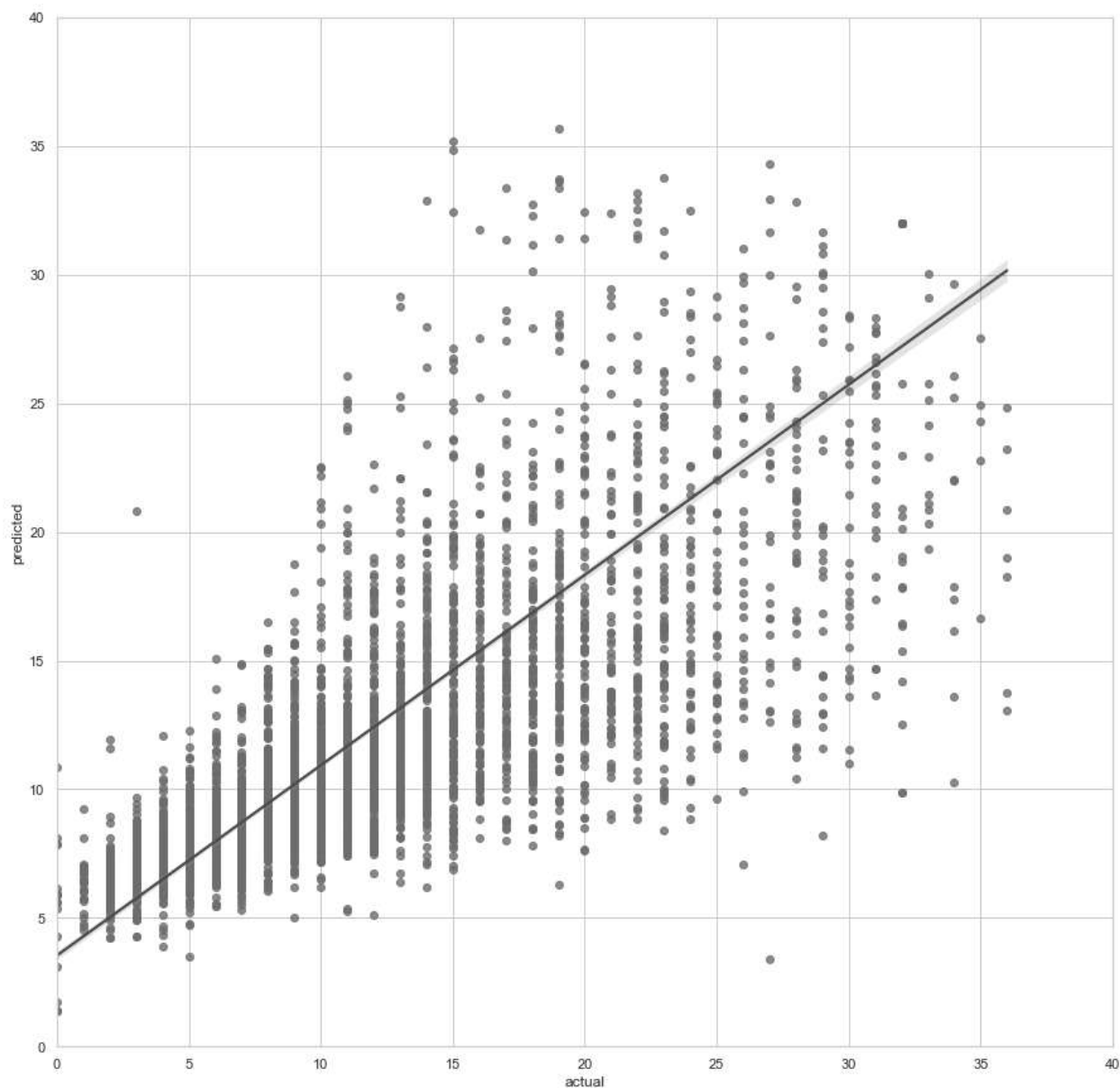
```
In [36]: 1 # Create a `results` dataframe
2
3 results = pd.DataFrame(data={"actual": Y_test["trip_duration"],
4                               "predicted": Y_pred_test.ravel()})
5 results["residual"] = results["actual"] - results["predicted"]
6 results.head()
```

Out[36]:

	actual	predicted	residual
<b>5818</b>	18.0	15.675123	2.324877
<b>18134</b>	32.0	32.000118	-0.000118
<b>4655</b>	6.0	7.328350	-1.328350
<b>7378</b>	16.0	18.358000	-2.358000
<b>13914</b>	11.0	11.550253	-0.550253

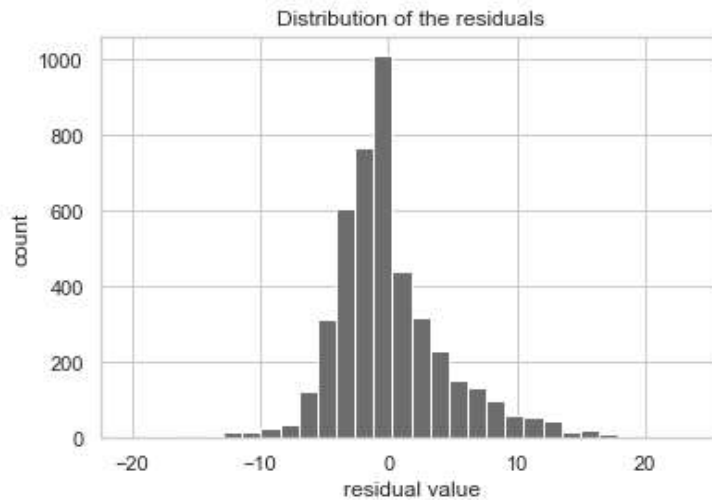
## Visualize model results

```
In [37]: 1 # Create a scatterplot to visualize `predicted` over `actual`  
2  
3 sns.set(style='whitegrid')  
4 f = plt.figure()  
5 f.set_figwidth(15)  
6 f.set_figheight(15)  
7 sns.regplot(x="actual",  
8             y="predicted",  
9             data=results, line_kws={"color": "red"})  
10 plt.ylim(0, 40)  
11 plt.xlim(0,40)  
12 plt.show()
```

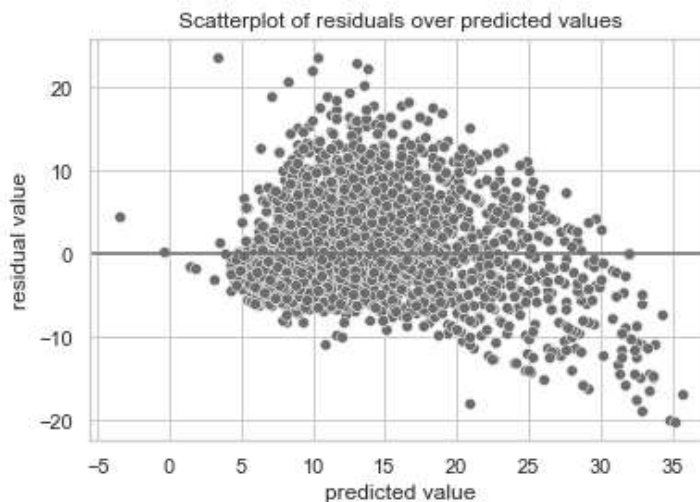




```
In [38]: 1 # Visualize the distribution of the `residuals`  
2  
3 plt.hist(results["residual"], bins=30)  
4 plt.title("Distribution of the residuals")  
5 plt.xlabel("residual value")  
6 plt.ylabel("count")  
7 plt.show()
```



```
In [39]: 1 # Create a scatterplot of `residuals` over `predicted`  
2  
3 sns.scatterplot(x="predicted", y="residual", data=results)  
4 plt.axhline(0)  
5 plt.title("Scatterplot of residuals over predicted values")  
6 plt.xlabel("predicted value")  
7 plt.ylabel("residual value")  
8 plt.show()
```



## Conclusion

Multiple linear regression is a powerful tool to estimate a dependent continuous variable from several independent variables.

Exploratory data analysis is useful for selecting both numeric and categorical features for multiple linear regression.

Fitting multiple linear regression models may require trial and error to select variables that fit an accurate model while maintaining model assumptions.

In [ ]:

1