

Business Case: Delhivery - Feature Engineering



Delhivery is the largest and fastest-growing fully integrated player in India by revenue in Fiscal 2021. They aim to build the operating system for commerce, through a combination of world-class infrastructure, logistics operations of the highest quality, and cutting-edge engineering and technology capabilities.

```
In [4]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
```

```
In [5]: data = pd.read_csv('delhivery_data.csv')
```

```
In [6]: data.head()
```

Out[6]:

	data	trip_creation_time	route_schedule_uuid	route_type	trip_uuid	source_center	source_name	destination_center
0	training	2018-09-20 02:35:36.476840	thanos::sroute:eb7bfc78-b351-4c0e-a951-fa3d5c3...	Carting	trip-153741093647649320	IND388121AAA	Anand_VUNagar_DC (Gujarat)	IND388620/
1	training	2018-09-20 02:35:36.476840	thanos::sroute:eb7bfc78-b351-4c0e-a951-fa3d5c3...	Carting	trip-153741093647649320	IND388121AAA	Anand_VUNagar_DC (Gujarat)	IND388620/
2	training	2018-09-20 02:35:36.476840	thanos::sroute:eb7bfc78-b351-4c0e-a951-fa3d5c3...	Carting	trip-153741093647649320	IND388121AAA	Anand_VUNagar_DC (Gujarat)	IND388620/
3	training	2018-09-20 02:35:36.476840	thanos::sroute:eb7bfc78-b351-4c0e-a951-fa3d5c3...	Carting	trip-153741093647649320	IND388121AAA	Anand_VUNagar_DC (Gujarat)	IND388620/
4	training	2018-09-20 02:35:36.476840	thanos::sroute:eb7bfc78-b351-4c0e-a951-fa3d5c3...	Carting	trip-153741093647649320	IND388121AAA	Anand_VUNagar_DC (Gujarat)	IND388620/

5 rows × 24 columns



Features of the dataset:

Field	Description
data	Indicates whether the data is testing or training data
trip_creation_time	Timestamp of trip creation
route_schedule_uuid	Unique ID for a particular route schedule
route_type	Transportation type
FTL	Full Truck Load: FTL shipments reach the destination sooner with no other pickups or drop-offs

Field	Description
Carting	Handling system consisting of small vehicles (carts)
trip_uuid	Unique ID for a particular trip (may include different source and destination centers)
source_center	Source ID of trip origin
source_name	Source Name of trip origin
destination_center	Destination ID
destination_name	Destination Name
od_start_time	Trip start time
od_end_time	Trip end time
start_scan_to_end_scan	Time taken to deliver from source to destination
is_cutoff	Unknown field
cutoff_factor	Unknown field
cutoff_timestamp	Unknown field
actual_distance_to_destination	Distance in Kms between source and destination warehouse
actual_time	Actual time taken to complete the delivery (Cumulative)
osrm_time	Open-source routing engine time calculator (Cumulative)
osrm_distance	Open-source routing engine computed distance (Cumulative)
factor	Unknown field
segment_actual_time	Segment time. Time taken by the subset of the package delivery
segment_osrm_time	OSRM segment time. Time taken by the subset of the package delivery
segment_osrm_distance	OSRM distance. Distance covered by the subset of the package delivery
segment_factor	Unknown field

Key Concepts

1. Feature Creation

- Definition: The process of creating new features from existing data to improve model performance.
- Importance: Helps in enhancing the predictive power of the model by providing more relevant information.

2. Relationship between Features

- Definition: Understanding how features interact and influence each other.
- Importance: Can help in feature selection and engineering to improve model accuracy.

3. Column Normalization / Column Standardization

- Definition: Techniques used to scale features to a similar range.
- Importance: Essential for algorithms that rely on the distance between data points, like KNN.

4. Handling Categorical Values

- Definition: Methods to convert categorical data into a numerical format.
- Importance: Necessary for machine learning algorithms that require numerical input.

5. Missing Values and Outlier Treatment

- Definition: Strategies to handle missing data and identify outliers.
- Types of Outliers:
 - **Univariate Outliers:** Extreme values in a single feature.
 - **Multivariate Outliers:** Unusual combinations of values across multiple features.
- Importance: Ensures data quality and improves model robustness.

Exploratory Data Analysis

```
In [7]: data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 144867 entries, 0 to 144866
Data columns (total 24 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   data              144867 non-null   object  
 1   trip_creation_time 144867 non-null   object  
 2   route_schedule_uuid 144867 non-null   object  
 3   route_type          144867 non-null   object  
 4   trip_uuid           144867 non-null   object  
 5   source_center        144867 non-null   object  
 6   source_name          144574 non-null   object  
 7   destination_center   144867 non-null   object  
 8   destination_name     144606 non-null   object  
 9   od_start_time        144867 non-null   object  
 10  od_end_time         144867 non-null   object  
 11  start_scan_to_end_scan 144867 non-null   float64
 12  is_cutoff            144867 non-null   bool    
 13  cutoff_factor         144867 non-null   int64  
 14  cutoff_timestamp      144867 non-null   object  
 15  actual_distance_to_destination 144867 non-null   float64
 16  actual_time           144867 non-null   float64
 17  osrm_time             144867 non-null   float64
 18  osrm_distance          144867 non-null   float64
 19  factor                144867 non-null   float64
 20  segment_actual_time    144867 non-null   float64
 21  segment_osrm_time      144867 non-null   float64
 22  segment_osrm_distance 144867 non-null   float64
 23  segment_factor          144867 non-null   float64
dtypes: bool(1), float64(10), int64(1), object(12)
memory usage: 25.6+ MB
```

```
In [9]: data.shape
```

```
Out[9]: (144867, 24)
```

```
In [10]: print(f"TOTAL ROWS : {data.shape[0]}")
print(f"TOTAL COLUMNS : {data.shape[1]}")
```

TOTAL ROWS : 144867

TOTAL COLUMNS : 24

In [11]: `print(f"SIZE OF DataFrame : {data.size}")`

SIZE OF DataFrame : 3476808

In [12]: `print(f"Index of the DataFrame : {data.index}")`

Index of the DataFrame : RangeIndex(start=0, stop=144867, step=1)

In [14]: `print(f"Columns : {data.columns}")`

```
Columns : Index(['data', 'trip_creation_time', 'route_schedule_uuid', 'route_type',
    'trip_uuid', 'source_center', 'source_name', 'destination_center',
    'destination_name', 'od_start_time', 'od_end_time',
    'start_scan_to_end_scan', 'is_cutoff', 'cutoff_factor',
    'cutoff_timestamp', 'actual_distance_to_destination', 'actual_time',
    'osrm_time', 'osrm_distance', 'factor', 'segment_actual_time',
    'segment_osrm_time', 'segment_osrm_distance', 'segment_factor'],
   dtype='object')
```

In [15]: `data.describe()`

	start_scan_to_end_scan	cutoff_factor	actual_distance_to_destination	actual_time	osrm_time	osrm_distance	factor
count	144867.000000	144867.000000	144867.000000	144867.000000	144867.000000	144867.000000	144867.000000
mean	961.262986	232.926567	234.073372	416.927527	213.868272	284.771297	2.120107
std	1037.012769	344.755577	344.990009	598.103621	308.011085	421.119294	1.715421
min	20.000000	9.000000	9.000045	9.000000	6.000000	9.008200	0.144000
25%	161.000000	22.000000	23.355874	51.000000	27.000000	29.914700	1.604264
50%	449.000000	66.000000	66.126571	132.000000	64.000000	78.525800	1.857143
75%	1634.000000	286.000000	286.708875	513.000000	257.000000	343.193250	2.213483
max	7898.000000	1927.000000	1927.447705	4532.000000	1686.000000	2326.199100	77.387097



```
In [16]: data.describe().T
```

Out[16]:

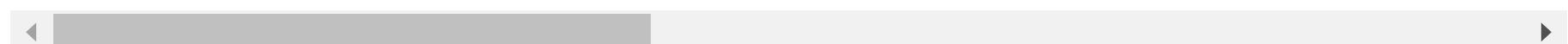
	count	mean	std	min	25%	50%	75%	max
start_scan_to_end_scan	144867.0	961.262986	1037.012769	20.000000	161.000000	449.000000	1634.000000	7898.000000
cutoff_factor	144867.0	232.926567	344.755577	9.000000	22.000000	66.000000	286.000000	1927.000000
actual_distance_to_destination	144867.0	234.073372	344.990009	9.000045	23.355874	66.126571	286.708875	1927.447705
actual_time	144867.0	416.927527	598.103621	9.000000	51.000000	132.000000	513.000000	4532.000000
osrm_time	144867.0	213.868272	308.011085	6.000000	27.000000	64.000000	257.000000	1686.000000
osrm_distance	144867.0	284.771297	421.119294	9.008200	29.914700	78.525800	343.193250	2326.199100
factor	144867.0	2.120107	1.715421	0.144000	1.604264	1.857143	2.213483	77.387097
segment_actual_time	144867.0	36.196111	53.571158	-244.000000	20.000000	29.000000	40.000000	3051.000000
segment_osrm_time	144867.0	18.507548	14.775960	0.000000	11.000000	17.000000	22.000000	1611.000000
segment_osrm_distance	144867.0	22.829020	17.860660	0.000000	12.070100	23.513000	27.813250	2191.403700
segment_factor	144867.0	2.218368	4.847530	-23.444444	1.347826	1.684211	2.250000	574.250000

```
In [17]: data.isnull()
```

Out[17]:

	data	trip_creation_time	route_schedule_uuid	route_type	trip_uuid	source_center	source_name	destination_center	destination_i
0	False	False	False	False	False	False	False	False	False
1	False	False	False	False	False	False	False	False	False
2	False	False	False	False	False	False	False	False	False
3	False	False	False	False	False	False	False	False	False
4	False	False	False	False	False	False	False	False	False
...
144862	False	False	False	False	False	False	False	False	False
144863	False	False	False	False	False	False	False	False	False
144864	False	False	False	False	False	False	False	False	False
144865	False	False	False	False	False	False	False	False	False
144866	False	False	False	False	False	False	False	False	False

144867 rows × 24 columns



In [18]: `data.isnull().sum()`

```
Out[18]: data          0
         trip_creation_time      0
         route_schedule_uuid      0
         route_type              0
         trip_uuid                0
         source_center            0
         source_name               293
         destination_center        0
         destination_name           261
         od_start_time             0
         od_end_time               0
         start_scan_to_end_scan     0
         is_cutoff                 0
         cutoff_factor              0
         cutoff_timestamp            0
         actual_distance_to_destination 0
         actual_time                  0
         osrm_time                   0
         osrm_distance                0
         factor                      0
         segment_actual_time          0
         segment_osrm_time             0
         segment_osrm_distance          0
         segment_factor                0
         dtype: int64
```

```
In [21]: # Note: There are some data.duplicated()
```

```
In [22]: data.duplicated()
```

```
Out[22]: 0      False
         1      False
         2      False
         3      False
         4      False
        ...
144862  False
144863  False
144864  False
144865  False
144866  False
Length: 144867, dtype: bool
```

```
In [23]: data.duplicated().sum()
```

```
Out[23]: 0
```

```
In [24]: # Note: There are zero duplicate values in the dataset
```

```
In [26]: #Taking a copy of data into data_copy
data_copy = data.copy()
```

```
In [27]: data.nunique()
```

```
Out[27]: data          2
trip_creation_time      14817
route_schedule_uuid     1504
route_type              2
trip_uuid                14817
source_center            1508
source_name               1498
destination_center       1481
destination_name         1468
od_start_time            26369
od_end_time               26369
start_scan_to_end_scan   1915
is_cutoff                 2
cutoff_factor             501
cutoff_timestamp          93180
actual_distance_to_destination 144515
actual_time                  3182
osrm_time                   1531
osrm_distance                138046
factor                      45641
segment_actual_time        747
segment_osrm_time           214
segment_osrm_distance      113799
segment_factor                5675
dtype: int64
```

```
In [28]: from scipy import stats
import warnings

# Ignore warnings
warnings.filterwarnings('ignore')
```

```
In [29]: #Dropping unknown fields
unknown_fields = ['is_cutoff', 'cutoff_factor', 'cutoff_timestamp', 'factor', 'segment_factor']
data = data.drop(columns = unknown_fields)
```

```
In [31]: #How many unique entries present in each column ?
for i in data.columns:
    print(f"Unique entries for column {i[:30]} = {data[i].nunique()}")
```

```
Unique entries for column data = 2
Unique entries for column trip_creation_time = 14817
Unique entries for column route_schedule_uuid = 1504
Unique entries for column route_type = 2
Unique entries for column trip_uuid = 14817
Unique entries for column source_center = 1508
Unique entries for column source_name = 1498
Unique entries for column destination_center = 1481
Unique entries for column destination_name = 1468
Unique entries for column od_start_time = 26369
Unique entries for column od_end_time = 26369
Unique entries for column start_scan_to_end_scan = 1915
Unique entries for column actual_distance_to_destination = 144515
Unique entries for column actual_time = 3182
Unique entries for column osrm_time = 1531
Unique entries for column osrm_distance = 138046
Unique entries for column segment_actual_time = 747
Unique entries for column segment_osrm_time = 214
Unique entries for column segment_osrm_distance = 113799
```

```
In [33]: #For all those columns where number of unique entries is 2, converting the datatype of columns to category
data['data'] = data['data'].astype('category')
data['route_type'] = data['route_type'].astype('category')
```

```
In [34]: floating_columns = ['actual_distance_to_destination', 'actual_time', 'osrm_time', 'osrm_distance',
                           'segment_actual_time', 'segment_osrm_time', 'segment_osrm_distance']
for i in floating_columns:
    print(i, data[i].max())
```

```
actual_distance_to_destination 1927.4477046975032
actual_time 4532.0
osrm_time 1686.0
osrm_distance 2326.1991000000003
segment_actual_time 3051.0
segment_osrm_time 1611.0
segment_osrm_distance 2191.4037000000003
```

```
In [35]: #We can update the datatype to float32 since the maximum value entry is small
for i in floating_columns:
    data[i] = data[i].astype('float32')
```

```
In [36]: #Updating the datatype of the datetime columns
datetime_columns = ['trip_creation_time', 'od_start_time', 'od_end_time']
for i in datetime_columns:
    data[i] = pd.to_datetime(data[i])
```

```
In [37]: data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 144867 entries, 0 to 144866
Data columns (total 19 columns):
 #   Column           Non-Null Count  Dtype  
 ---  -- 
 0   data              144867 non-null   category
 1   trip_creation_time 144867 non-null   datetime64[ns]
 2   route_schedule_uuid 144867 non-null   object  
 3   route_type          144867 non-null   category
 4   trip_uuid           144867 non-null   object  
 5   source_center        144867 non-null   object  
 6   source_name          144574 non-null   object  
 7   destination_center   144867 non-null   object  
 8   destination_name     144606 non-null   object  
 9   od_start_time        144867 non-null   datetime64[ns]
 10  od_end_time         144867 non-null   datetime64[ns]
 11  start_scan_to_end_scan 144867 non-null   float64
 12  actual_distance_to_destination 144867 non-null   float32
 13  actual_time          144867 non-null   float32
 14  osrm_time            144867 non-null   float32
 15  osrm_distance        144867 non-null   float32
 16  segment_actual_time  144867 non-null   float32
 17  segment_osrm_time    144867 non-null   float32
 18  segment_osrm_distance 144867 non-null   float32
dtypes: category(2), datetime64[ns](3), float32(7), float64(1), object(6)
memory usage: 15.2+ MB
```

```
In [39]: #memory usage reduced from 25.6+ MB to 15.2+ MB
```

```
In [40]: #What is the time period for which the data is given ?
data['trip_creation_time'].min(), data['od_end_time'].max()
```

```
Out[40]: (Timestamp('2018-09-12 00:00:16.535741'),
           Timestamp('2018-10-08 03:00:24.353479'))
```

```
In [41]: data.isnull().sum()
```

```
Out[41]: data          0
         trip_creation_time 0
         route_schedule_uuid 0
         route_type          0
         trip_uuid            0
         source_center         0
         source_name           293
         destination_center    0
         destination_name      261
         od_start_time         0
         od_end_time           0
         start_scan_to_end_scan 0
         actual_distance_to_destination 0
         actual_time            0
         osrm_time              0
         osrm_distance          0
         segment_actual_time    0
         segment_osrm_time      0
         segment_osrm_distance   0
         dtype: int64
```

```
In [42]: missing_source_name = data.loc[data['source_name'].isnull(), 'source_center'].unique()
missing_source_name
```

```
Out[42]: array(['IND342902A1B', 'IND577116AAA', 'IND282002AAD', 'IND465333A1B',
               'IND841301AAC', 'IND509103AAC', 'IND126116AAA', 'IND331022A1B',
               'IND505326AAB', 'IND852118A1B'], dtype=object)
```

```
In [43]: for i in missing_source_name:
    unique_source_name = data.loc[data['source_center'] == i, 'source_name'].unique()
    if pd.isna(unique_source_name):
        print("Source Center :", i, "-" * 10, "Source Name :", 'Not Found')
    else :
        print("Source Center :", i, "-" * 10, "Source Name :", unique_source_name)
```

```
Source Center : IND342902A1B ----- Source Name : Not Found
Source Center : IND577116AAA ----- Source Name : Not Found
Source Center : IND282002AAD ----- Source Name : Not Found
Source Center : IND465333A1B ----- Source Name : Not Found
Source Center : IND841301AAC ----- Source Name : Not Found
Source Center : IND509103AAC ----- Source Name : Not Found
Source Center : IND126116AAA ----- Source Name : Not Found
Source Center : IND331022A1B ----- Source Name : Not Found
Source Center : IND505326AAB ----- Source Name : Not Found
Source Center : IND852118A1B ----- Source Name : Not Found
```

```
In [45]: for i in missing_source_name:
    unique_destination_name = data.loc[data['destination_center'] == i, 'destination_name'].unique()
    if (pd.isna(unique_source_name)) or (unique_source_name.size == 0):
        print("Destination Center :", i, "-" * 10, "Destination Name :", 'Not Found')
    else :
        print("Destination Center :", i, "-" * 10, "Destination Name :", unique_destination_name)
```

```
Destination Center : IND342902A1B ----- Destination Name : Not Found
Destination Center : IND577116AAA ----- Destination Name : Not Found
Destination Center : IND282002AAD ----- Destination Name : Not Found
Destination Center : IND465333A1B ----- Destination Name : Not Found
Destination Center : IND841301AAC ----- Destination Name : Not Found
Destination Center : IND509103AAC ----- Destination Name : Not Found
Destination Center : IND126116AAA ----- Destination Name : Not Found
Destination Center : IND331022A1B ----- Destination Name : Not Found
Destination Center : IND505326AAB ----- Destination Name : Not Found
Destination Center : IND852118A1B ----- Destination Name : Not Found
```

```
In [47]: missing_destination_name = data.loc[data['destination_name'].isnull(), 'destination_center'].unique()
missing_destination_name
```

```
Out[47]: array(['IND342902A1B', 'IND577116AAA', 'IND282002AAD', 'IND465333A1B',
   'IND841301AAC', 'IND505326AAB', 'IND852118A1B', 'IND126116AAA',
   'IND509103AAC', 'IND221005A1A', 'IND250002AAC', 'IND331001A1C',
   'IND122015AAC'], dtype=object)
```

```
In [48]: #The IDs for which the source name is missing, are all those IDs for destination also missing ?
np.all(data.loc[data['source_name'].isnull(), 'source_center'].isin(missing_destination_name))
```

```
Out[48]: False
```

```
In [50]: #Treating missing destination names and source names
count = 1
for i in missing_destination_name:
    data.loc[data['destination_center'] == i, 'destination_name'] = data.loc[data['destination_center'] == i, 'destination_name'].unique()
    count += 1
```

```
In [51]: d = {}
for i in missing_source_name:
    d[i] = data.loc[data['source_center'] == i, 'source_name'].unique()
for idx, val in d.items():
    if len(val) == 0:
        d[idx] = [f'location_{count}']
        count += 1
d2 = {}
for idx, val in d.items():
    d2[idx] = val[0]
for i, v in d2.items():
    print(i, v)
```

```
IND342902A1B location_1
IND577116AAA location_2
IND282002AAD location_3
IND465333A1B location_4
IND841301AAC location_5
IND509103AAC location_9
IND126116AAA location_8
IND331022A1B location_14
IND505326AAB location_6
IND852118A1B location_7
```

```
In [52]: for i in missing_source_name:
    data.loc[data['source_center'] == i, 'source_name'] = data.loc[data['source_center'] == i, 'source_name'].replace(np.nan,
```

```
In [54]: data.isna().sum()
```

```
Out[54]: data          0  
trip_creation_time      0  
route_schedule_uuid      0  
route_type                0  
trip_uuid                  0  
source_center              0  
source_name                  0  
destination_center          0  
destination_name            0  
od_start_time               0  
od_end_time                 0  
start_scan_to_end_scan      0  
actual_distance_to_destination 0  
actual_time                  0  
osrm_time                      0  
osrm_distance                  0  
segment_actual_time           0  
segment_osrm_time             0  
segment_osrm_distance         0  
dtype: int64
```

```
In [55]: data.describe()
```

Out[55]:

	trip_creation_time	od_start_time	od_end_time	start_scan_to_end_scan	actual_distance_to_destination	actual_time	
count	144867	144867	144867	144867.000000	144867.000000	144867.000000	144
mean	2018-09-22 13:34:23.659819264	2018-09-22 18:02:45.855230720	2018-09-23 10:04:31.395393024	961.262986	234.073380	416.927521	
min	2018-09-12 00:00:16.535741	2018-09-12 00:00:16.535741	2018-09-12 00:50:10.814399	20.000000	9.000046	9.000000	
25%	2018-09-17 03:20:51.775845888	2018-09-17 08:05:40.886155008	2018-09-18 01:48:06.410121984	161.000000	23.355875	51.000000	
50%	2018-09-22 04:24:27.932764928	2018-09-22 08:53:00.116656128	2018-09-23 03:13:03.520212992	449.000000	66.126572	132.000000	
75%	2018-09-27 17:57:56.350054912	2018-09-27 22:41:50.285857024	2018-09-28 12:49:06.054018048	1634.000000	286.708878	513.000000	
max	2018-10-03 23:59:42.701692	2018-10-06 04:27:23.392375	2018-10-08 03:00:24.353479	7898.000000	1927.447754	4532.000000	1
std	Nan	Nan	Nan	1037.012769	344.990021	598.103638	

◀ ▶

In [56]: `data.describe(include = 'object')`

Out[56]:

	route_schedule_uuid	trip_uuid	source_center	source_name	destination_center	destination_name	
count	144867	144867	144867	144867	144867	144867	144867
unique	1504	14817	1508	1508	1481	1481	
top	thanos::sroute:4029a8a2-6c74-4b7e-a6d8-f9e069f...	trip-153811219535896559	IND000000ACB	Gurgaon_Bilaspur_HB (Haryana)	IND000000ACB	Gurgaon_Bilaspur_HB (Haryana)	
freq	1812	101	23347	23347	15192	15192	

Merging of rows and aggregation of fields

How to begin"

Since delivery details of one package are divided into several rows (think of it as connecting flights to reach a particular destination). Now think about how we should treat their fields if we combine these rows? What aggregation would make sense if we merge. What would happen to the numeric fields if we merge the rows.

```
In [58]: grouping_1 = ['trip_uuid', 'source_center', 'destination_center']
data_1 = data.groupby(by = grouping_1, as_index = False).agg({'data' : 'first',
                                                               'route_type' : 'first',
                                                               'trip_creation_time' : 'first',
                                                               'source_name' : 'first',
                                                               'destination_name' : 'last',
                                                               'od_start_time' : 'first',
                                                               'od_end_time' : 'first',
                                                               'start_scan_to_end_scan' : 'first',
                                                               'actual_distance_to_destination' : 'last',
                                                               'actual_time' : 'last',
                                                               'osrm_time' : 'last',
                                                               'osrm_distance' : 'last',
                                                               'segment_actual_time' : 'sum',
                                                               'segment_osrm_time' : 'sum',
                                                               'segment_osrm_distance' : 'sum'})
```

data_1

Out[58]:

	trip_uuid	source_center	destination_center	data	route_type	trip_creation_time	source_name	des
0	trip-153671041653548748	IND209304AAA	IND000000ACB	training	FTL	2018-09-12 00:00:16.535741	Kanpur_Central_H_6 (Uttar Pradesh)	Gurgaon
1	trip-153671041653548748	IND462022AAA	IND209304AAA	training	FTL	2018-09-12 00:00:16.535741	Bhopal_Tnrsport_H (Madhya Pradesh)	Kanpur
2	trip-153671042288605164	IND561203AAB	IND562101AAA	training	Carting	2018-09-12 00:00:22.886430	Doddablpur_ChikaDPP_D (Karnataka)	Chikballapur
3	trip-153671042288605164	IND572101AAA	IND561203AAB	training	Carting	2018-09-12 00:00:22.886430	Tumkur_Veersagr_I (Karnataka)	Doddablpur
4	trip-153671043369099517	IND000000ACB	IND160002AAC	training	FTL	2018-09-12 00:00:33.691250	Gurgaon_Bilaspur_HB (Haryana)	Chandigarh
...
26363	trip-153861115439069069	IND628204AAA	IND627657AAA	test	Carting	2018-10-03 23:59:14.390954	Tirchchndr_Shnmgrpm_D (Tamil Nadu)	Thisayanvil
26364	trip-153861115439069069	IND628613AAA	IND627005AAA	test	Carting	2018-10-03 23:59:14.390954	Peikulam_SriVnktpm_D (Tamil Nadu)	Tirunelveli
26365	trip-153861115439069069	IND628801AAA	IND628204AAA	test	Carting	2018-10-03 23:59:14.390954	Eral_Busstand_D (Tamil Nadu)	Tirchchndi
26366	trip-153861118270144424	IND583119AAA	IND583101AAA	test	FTL	2018-10-03 23:59:42.701692	Sandur_WrdN1DPP_D (Karnataka)	Bellary
26367	trip-153861118270144424	IND583201AAA	IND583119AAA	test	FTL	2018-10-03 23:59:42.701692	Hospet (Karnataka)	Sandur

26368 rows × 18 columns



In [60]:

```
#Calculate the time taken between od_start_time and od_end_time and keep it as a feature. Drop the original columns, if required
data_1['od_total_time'] = data_1['od_end_time'] - data_1['od_start_time']
data_1.drop(columns = ['od_end_time', 'od_start_time'], inplace = True)
```

```
data_1['od_total_time'] = data_1['od_total_time'].apply(lambda x : round(x.total_seconds() / 60.0, 2))
data_1['od_total_time'].head()
```

```
Out[60]: 0    1260.60
1    999.51
2    58.83
3   122.78
4   834.64
Name: od_total_time, dtype: float64
```

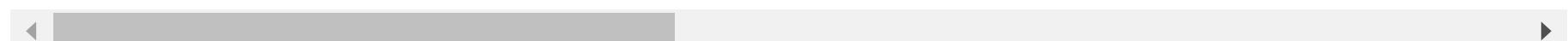
```
In [61]: data_2 = data_1.groupby(by = 'trip_uuid', as_index = False).agg({'source_center' : 'first',
                                                               'destination_center' : 'last',
                                                               'data' : 'first',
                                                               'route_type' : 'first',
                                                               'trip_creation_time' : 'first',
                                                               'source_name' : 'first',
                                                               'destination_name' : 'last',
                                                               'od_total_time' : 'sum',
                                                               'start_scan_to_end_scan' : 'sum',
                                                               'actual_distance_to_destination' : 'sum',
                                                               'actual_time' : 'sum',
                                                               'osrm_time' : 'sum',
                                                               'osrm_distance' : 'sum',
                                                               'segment_actual_time' : 'sum',
                                                               'segment_osrm_time' : 'sum',
                                                               'segment_osrm_distance' : 'sum'})
```

```
data_2
```

Out[61]:

	trip_uuid	source_center	destination_center	data	route_type	trip_creation_time	source_name	de
0	trip-153671041653548748	IND209304AAA	IND209304AAA	training	FTL	2018-09-12 00:00:16.535741	Kanpur_Central_H_6 (Uttar Pradesh)	Kan
1	trip-153671042288605164	IND561203AAB	IND561203AAB	training	Carting	2018-09-12 00:00:22.886430	Doddablpur_ChikaDPP_D (Karnataka)	Doddablpur
2	trip-153671043369099517	IND000000ACB	IND000000ACB	training	FTL	2018-09-12 00:00:33.691250	Gurgaon_Bilaspur_HB (Haryana)	Gurgaon
3	trip-153671046011330457	IND400072AAB	IND401104AAA	training	Carting	2018-09-12 00:01:00.113710	Mumbai Hub (Maharashtra)	Mumbai
4	trip-153671052974046625	IND583101AAA	IND583119AAA	training	FTL	2018-09-12 00:02:09.740725	Bellary_Dc (Karnataka)	Sandur
...
14812	trip-153861095625827784	IND160002AAC	IND160002AAC	test	Carting	2018-10-03 23:55:56.258533	Chandigarh_Mehmdpur_H (Punjab)	Chandigarh
14813	trip-153861104386292051	IND121004AAB	IND121004AAA	test	Carting	2018-10-03 23:57:23.863155	FBD_Balabgarh_DPC (Haryana)	Faridabad
14814	trip-153861106442901555	IND208006AAA	IND208006AAA	test	Carting	2018-10-03 23:57:44.429324	Kanpur_GovndNgr_DC (Uttar Pradesh)	Kanpur
14815	trip-153861115439069069	IND627005AAA	IND628204AAA	test	Carting	2018-10-03 23:59:14.390954	Tirunelveli_VdkkuSrt_I (Tamil Nadu)	Tiruchchirappalli
14816	trip-153861118270144424	IND583119AAA	IND583119AAA	test	FTL	2018-10-03 23:59:42.701692	Sandur_WrdN1DPP_D (Karnataka)	Sandur

14817 rows × 17 columns



Build some features to prepare the data for actual analysis. Extract features from the below fields

```
In [62]: #Source Name: Split and extract features out of destination. City-place-code (State)
def location_name_to_state(x):
    l = x.split('(')
    if len(l) == 1:
        return l[0]
    else:
        return l[1].replace(')', "")
```

```
In [63]: def location_name_to_city(x):
    if 'location' in x:
        return 'unknown_city'
    else:
        l = x.split()[0].split('_')
        if 'CCU' in x:
            return 'Kolkata'
        elif 'MAA' in x.upper():
            return 'Chennai'
        elif ('HBR' in x.upper()) or ('BLR' in x.upper()):
            return 'Bengaluru'
        elif 'FBD' in x.upper():
            return 'Faridabad'
        elif 'BOM' in x.upper():
            return 'Mumbai'
        elif 'DEL' in x.upper():
            return 'Delhi'
        elif 'OK' in x.upper():
            return 'Delhi'
        elif 'GZB' in x.upper():
            return 'Ghaziabad'
        elif 'GGN' in x.upper():
            return 'Gurgaon'
        elif 'AMD' in x.upper():
            return 'Ahmedabad'
        elif 'CJB' in x.upper():
            return 'Coimbatore'
        elif 'HYD' in x.upper():
            return 'Hyderabad'
        return l[0]
```

```
In [64]: def location_name_to_place(x):
    if 'location' in x:
        return x
    elif 'HBR' in x:
        return 'HBR Layout PC'
    else:
        l = x.split()[0].split('_', 1)
        if len(l) == 1:
            return 'unknown_place'
        else:
            return l[1]
```

```
In [65]: data_2['source_state'] = data_2['source_name'].apply(location_name_to_state)
data_2['source_state'].unique()
```

```
Out[65]: array(['Uttar Pradesh', 'Karnataka', 'Haryana', 'Maharashtra',
       'Tamil Nadu', 'Gujarat', 'Delhi', 'Telangana', 'Rajasthan',
       'Assam', 'Madhya Pradesh', 'West Bengal', 'Andhra Pradesh',
       'Punjab', 'Chandigarh', 'Goa', 'Jharkhand', 'Pondicherry',
       'Orissa', 'Uttarakhand', 'Himachal Pradesh', 'Kerala',
       'Arunachal Pradesh', 'Bihar', 'Chhattisgarh',
       'Dadra and Nagar Haveli', 'Jammu & Kashmir', 'Mizoram', 'Nagaland',
       'location_9', 'location_3', 'location_2', 'location_14',
       'location_7'], dtype=object)
```

```
In [66]: data_2['source_city'] = data_2['source_name'].apply(location_name_to_city)
print('No of source cities :', data_2['source_city'].nunique())
data_2['source_city'].unique()[:100]
```

```
No of source cities : 690
```

```
Out[66]: array(['Kanpur', 'Doddablpur', 'Gurgaon', 'Mumbai', 'Bellary', 'Chennai',
   'Bengaluru', 'Surat', 'Delhi', 'Pune', 'Faridabad', 'Shirala',
   'Hyderabad', 'Thirumalagiri', 'Gulbarga', 'Jaipur', 'Allahabad',
   'Guwahati', 'Narsinghpur', 'Shrirampur', 'Madakasira', 'Sonari',
   'Dindigul', 'Jalandhar', 'Chandigarh', 'Deoli', 'Pandharpur',
   'Kolkata', 'Bhandara', 'Kurnool', 'Bhiwandi', 'Bhatinda',
   'RoopNagar', 'Bantwal', 'Lalru', 'Kadi', 'Shahdol', 'Gangakher',
   'Durgapur', 'Vapi', 'Jamjodhpur', 'Jetpur', 'Mehsana', 'Jabalpur',
   'Junagadh', 'Gundlupet', 'Mysore', 'Goa', 'Bhopal', 'Sonipat',
   'Himmatnagar', 'Jamshedpur', 'Pondicherry', 'Anand', 'Udgir',
   'Nadiad', 'Villupuram', 'Purulia', 'Bhubaneshwar', 'Bamangola',
   'Tiruppattur', 'Kotdwara', 'Medak', 'Bangalore', 'Dhrangadhra',
   'Hospet', 'Ghumarwin', 'Agra', 'Sitapur', 'Canacona', 'Bilimora',
   'SultnBthry', 'Lucknow', 'Vellore', 'Bhuj', 'Dinhata',
   'Margherita', 'Boisar', 'Vizag', 'Tezpur', 'Koduru', 'Tirupati',
   'Pen', 'Ahmedabad', 'Faizabad', 'Gandhinagar', 'Anantapur',
   'Betul', 'Panskura', 'Rasipurm', 'Sankari', 'Jorhat', 'PNQ',
   'Srikakulam', 'Dehradun', 'Jassur', 'Sawantwadi', 'Shajapur',
   'Ludhiana', 'GreaterThane'], dtype=object)
```

```
In [67]: data_2['source_place'] = data_2['source_name'].apply(location_name_to_place)
data_2['source_place'].unique()[:100]
```

```
Out[67]: array(['Central_H_6', 'ChikaDPP_D', 'Bilaspur_HB', 'unknown_place', 'Dc',
       'Poonamallee', 'Chrompet_DPC', 'HBR Layout PC', 'Central_D_12',
       'Lajpat_IP', 'North_D_3', 'Balabgarh_DPC', 'Central_DPP_3',
       'Shamshbd_H', 'Xroad_D', 'Nehrugnj_I', 'Central_I_7',
       'Central_H_1', 'Nangli_IP', 'North', 'KndliDPP_D', 'Central_D_9',
       'DavkharRd_D', 'Bandel_D', 'RTCStand_D', 'Central_DPP_1',
       'KGAirprt_HB', 'North_D_2', 'Central_D_1', 'DC', 'Mthurard_L',
       'Mullanpr_DC', 'Central_DPP_2', 'RajCmplx_D', 'Belaghata_DPC',
       'RjnaiDPP_D', 'AbbasNgr_I', 'Mankoli_HB', 'DPC', 'Airport_H',
       'Hub', 'Gateway_HB', 'Tathawde_H', 'ChotiHvl_DC', 'Trmltmp_D',
       'OnkarDPP_D', 'Mehmdpur_H', 'KaranNGR_D', 'Sohagpur_D',
       'Chrompet_L', 'Busstand_D', 'Central_I_1', 'IndEstat_I', 'Court_D',
       'Panchot_IP', 'Adhartal_IP', 'DumDum_DPC', 'Bomsndra_HB',
       'Swamylyt_D', 'Yadvgiri_IP', 'Old', 'Kundli_H', 'Central_I_3',
       'Vasanthm_I', 'Poonamallee_HB', 'VUNagar_DC', 'NlgaonRd_D',
       'Bnnrghtha_L', 'Thirumtr_IP', 'GariDPP_D', 'Jogshwri_I',
       'KoilStrt_D', 'CotnGren_M', 'Nzbadrd_D', 'Dwaraka_D', 'Nelmngla_H',
       'NvygRDPP_D', 'Gndhichk_D', 'Central_D_3', 'Chowk_D', 'CharRsta_D',
       'Kollgpra_D', 'Peenya_IP', 'GndhiNgr_IP', 'Sanpada_I',
       'WrdN4DPP_D', 'Sakinaka_RP', 'CivilHPL_D', 'OstwlEmp_D',
       'Gajuwaka', 'Mhbhirab_D', 'MGRoad_D', 'Balajicly_I', 'BljiMrkt_D',
       'Dankuni_HB', 'Trnsport_H', 'Rakhial', 'Memnagar', 'East_I_21',
       'Mithakal_D'], dtype=object)
```

```
In [68]: #Destination Name: Split and extract features out of destination. City-place-code (State)
data_2['destination_state'] = data_2['destination_name'].apply(location_name_to_state)
data_2['destination_state'].head(10)
```

```
Out[68]: 0    Uttar Pradesh
1    Karnataka
2    Haryana
3    Maharashtra
4    Karnataka
5    Tamil Nadu
6    Tamil Nadu
7    Karnataka
8    Gujarat
9    Delhi
Name: destination_state, dtype: object
```

```
In [69]: data_2['destination_city'] = data_2['destination_name'].apply(location_name_to_city)
data_2['destination_city'].head()
```

```
Out[69]: 0      Kanpur
1    Doddablpur
2      Gurgaon
3      Mumbai
4      Sandur
Name: destination_city, dtype: object
```

```
In [70]: data_2['destination_place'] = data_2['destination_name'].apply(location_name_to_place)
data_2['destination_place'].head()
```

```
Out[70]: 0    Central_H_6
1    ChikaDPP_D
2    Bilaspur_HB
3    MiraRd_IP
4    WrdN1DPP_D
Name: destination_place, dtype: object
```

```
In [71]: #Trip_creation_time: Extract features Like month, year and day etc
data_2['trip_creation_date'] = pd.to_datetime(data_2['trip_creation_time'].dt.date)
data_2['trip_creation_date'].head()
```

```
Out[71]: 0    2018-09-12
1    2018-09-12
2    2018-09-12
3    2018-09-12
4    2018-09-12
Name: trip_creation_date, dtype: datetime64[ns]
```

```
In [72]: data_2['trip_creation_day'] = data_2['trip_creation_time'].dt.day
data_2['trip_creation_day'] = data_2['trip_creation_day'].astype('int8')
data_2['trip_creation_day'].head()
```

```
Out[72]: 0    12
         1    12
         2    12
         3    12
         4    12
Name: trip_creation_day, dtype: int8
```

```
In [73]: data_2['trip_creation_month'] = data_2['trip_creation_time'].dt.month
          data_2['trip_creation_month'] = data_2['trip_creation_month'].astype('int8')
          data_2['trip_creation_month'].head()
```

```
Out[73]: 0    9
         1    9
         2    9
         3    9
         4    9
Name: trip_creation_month, dtype: int8
```

```
In [74]: data_2['trip_creation_year'] = data_2['trip_creation_time'].dt.year
          data_2['trip_creation_year'] = data_2['trip_creation_year'].astype('int16')
          data_2['trip_creation_year'].head()
```

```
Out[74]: 0    2018
         1    2018
         2    2018
         3    2018
         4    2018
Name: trip_creation_year, dtype: int16
```

```
In [75]: data_2['trip_creation_week'] = data_2['trip_creation_time'].dt.isocalendar().week
          data_2['trip_creation_week'] = data_2['trip_creation_week'].astype('int8')
          data_2['trip_creation_week'].head()
```

```
Out[75]: 0    37
         1    37
         2    37
         3    37
         4    37
Name: trip_creation_week, dtype: int8
```

```
In [76]: data_2['trip_creation_hour'] = data_2['trip_creation_time'].dt.hour  
data_2['trip_creation_hour'] = data_2['trip_creation_hour'].astype('int8')  
data_2['trip_creation_hour'].head()
```

```
Out[76]: 0    0  
1    0  
2    0  
3    0  
4    0  
Name: trip_creation_hour, dtype: int8
```

```
In [78]: #Finding the structure of data after data cleaning  
data_2.shape
```

```
Out[78]: (14817, 29)
```

```
In [80]: data_2.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 14817 entries, 0 to 14816
Data columns (total 29 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   trip_uuid        14817 non-null   object  
 1   source_center    14817 non-null   object  
 2   destination_center 14817 non-null   object  
 3   data             14817 non-null   category
 4   route_type       14817 non-null   category
 5   trip_creation_time 14817 non-null   datetime64[ns]
 6   source_name      14817 non-null   object  
 7   destination_name 14817 non-null   object  
 8   od_total_time    14817 non-null   float64 
 9   start_scan_to_end_scan 14817 non-null   float64 
 10  actual_distance_to_destination 14817 non-null   float32 
 11  actual_time      14817 non-null   float32 
 12  osrm_time        14817 non-null   float32 
 13  osrm_distance    14817 non-null   float32 
 14  segment_actual_time 14817 non-null   float32 
 15  segment_osrm_time 14817 non-null   float32 
 16  segment_osrm_distance 14817 non-null   float32 
 17  source_state     14817 non-null   object  
 18  source_city      14817 non-null   object  
 19  source_place     14817 non-null   object  
 20  destination_state 14817 non-null   object  
 21  destination_city 14817 non-null   object  
 22  destination_place 14817 non-null   object  
 23  trip_creation_date 14817 non-null   datetime64[ns]
 24  trip_creation_day 14817 non-null   int8    
 25  trip_creation_month 14817 non-null   int8    
 26  trip_creation_year 14817 non-null   int16   
 27  trip_creation_week 14817 non-null   int8    
 28  trip_creation_hour 14817 non-null   int8    
dtypes: category(2), datetime64[ns](2), float32(7), float64(2), int16(1), int8(4), object(11)
memory usage: 2.2+ MB
```

```
In [81]: data_2.describe().T
```

Out[81]:

	count	mean	min	25%	50%	75%	
trip_creation_time	14817	2018-09-22 12:44:19.555167744	2018-09-12 00:00:16.535741	2018-09-17 02:51:25.129125888	2018-09-22 04:02:35.066945024	2018-09-27 19:37:41.898427904	23:59
od_total_time	14817.0	531.69763	23.46	149.93	280.77	638.2	
start_scan_to_end_scan	14817.0	530.810016	23.0	149.0	280.0	637.0	
actual_distance_to_destination	14817.0	164.477829	9.002461	22.837238	48.474072	164.583206	21
actual_time	14817.0	357.143768	9.0	67.0	149.0	370.0	
osrm_time	14817.0	161.384018	6.0	29.0	60.0	168.0	
osrm_distance	14817.0	204.344711	9.0729	30.819201	65.618805	208.475006	28
segment_actual_time	14817.0	353.892273	9.0	66.0	147.0	367.0	
segment_osrm_time	14817.0	180.949783	6.0	31.0	65.0	185.0	
segment_osrm_distance	14817.0	223.201157	9.0729	32.654499	70.154404	218.802399	35
trip_creation_date	14817	2018-09-21 23:46:58.627252736	2018-09-12 00:00:00	2018-09-17 00:00:00	2018-09-22 00:00:00	2018-09-27 00:00:00	2
trip_creation_day	14817.0	18.37079	1.0	14.0	19.0	25.0	
trip_creation_month	14817.0	9.120672	9.0	9.0	9.0	9.0	
trip_creation_year	14817.0	2018.0	2018.0	2018.0	2018.0	2018.0	
trip_creation_week	14817.0	38.295944	37.0	38.0	38.0	39.0	
trip_creation_hour	14817.0	12.449821	0.0	4.0	14.0	20.0	

In [82]: `data_2.describe(include = object).T`

Out[82]:

	count	unique	top	freq
trip_uuid	14817	14817	trip-153671041653548748	1
source_center	14817	938	IND000000ACB	1063
destination_center	14817	1042	IND000000ACB	821
source_name	14817	938	Gurgaon_Bilaspur_HB (Haryana)	1063
destination_name	14817	1042	Gurgaon_Bilaspur_HB (Haryana)	821
source_state	14817	34	Maharashtra	2714
source_city	14817	690	Mumbai	1442
source_place	14817	761	Bilaspur_HB	1063
destination_state	14817	39	Maharashtra	2561
destination_city	14817	806	Mumbai	1548
destination_place	14817	850	Bilaspur_HB	821

In [83]:

```
#I am interested to know how many trips are created on the hourly basis  
data_2['trip_creation_hour'].unique()
```

Out[83]:

```
array([ 0,  1,  2,  3,  4,  5,  6,  7,  8,  9, 10, 11, 12, 13, 14, 15, 16,  
       17, 18, 19, 20, 21, 22, 23], dtype=int8)
```

In [84]:

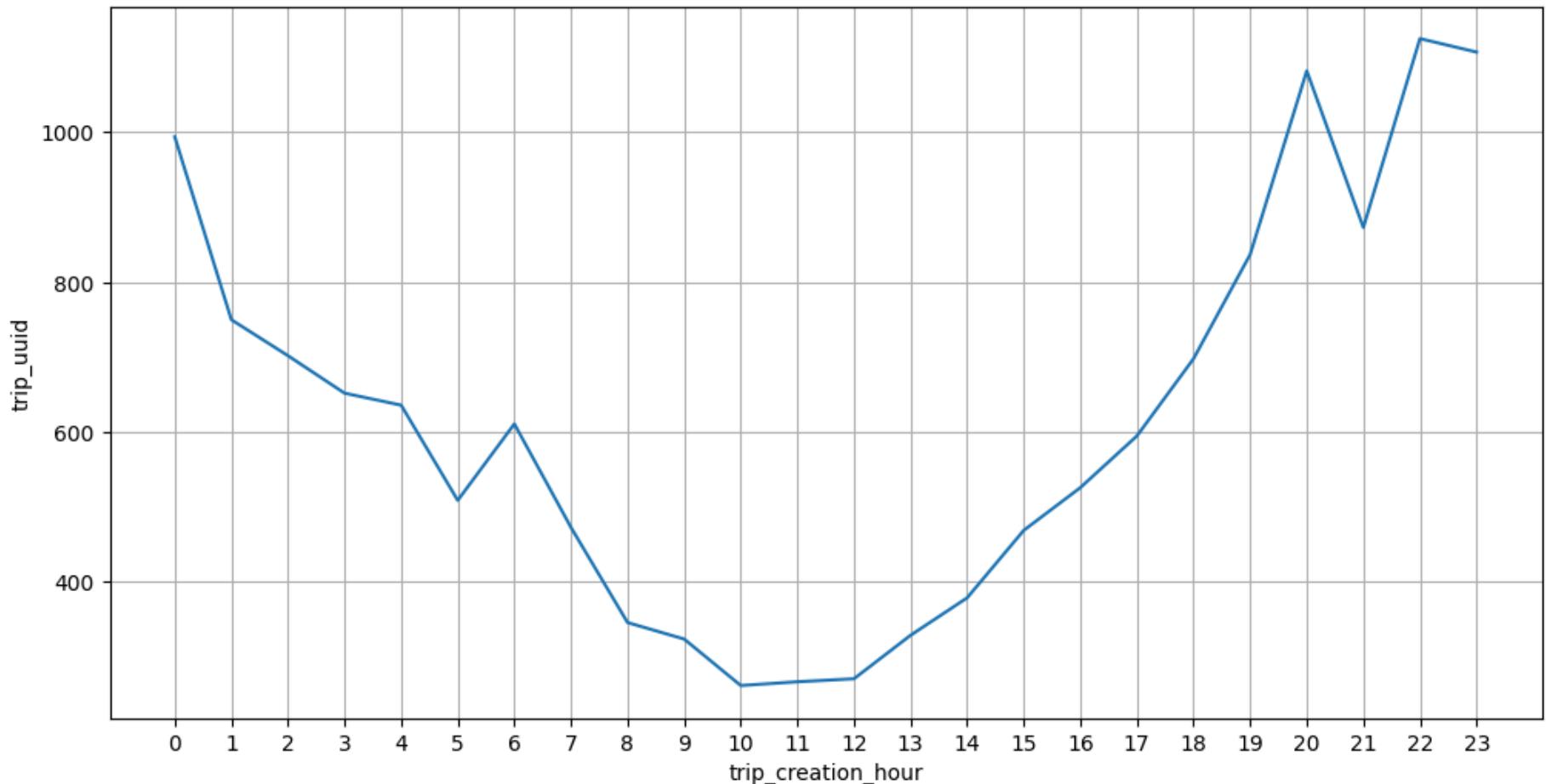
```
df_hour = data_2.groupby(by = 'trip_creation_hour')['trip_uuid'].count().to_frame().reset_index()  
df_hour.head()
```

```
Out[84]:    trip_creation_hour  trip_uuid
```

0	0	994
1	1	750
2	2	702
3	3	652
4	4	636

```
In [85]: plt.figure(figsize = (12, 6))
sns.lineplot(data = df_hour,
              x = df_hour['trip_creation_hour'],
              y = df_hour['trip_uuid'],
              markers = '*')
plt.xticks(np.arange(0,24))
plt.grid('both')
plt.plot()
```

```
Out[85]: []
```



```
In [86]: #It can be inferred from the above plot that the number of trips start increasing after the noon, becomes maximum at 10 P.M and then starts decreasing.
```

```
In [87]: #I am interested to know how many trips are created for different days of the month  
data_2['trip_creation_day'].unique()
```

```
Out[87]: array([12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28,  
   29, 30,  1,  2,  3], dtype=int8)
```

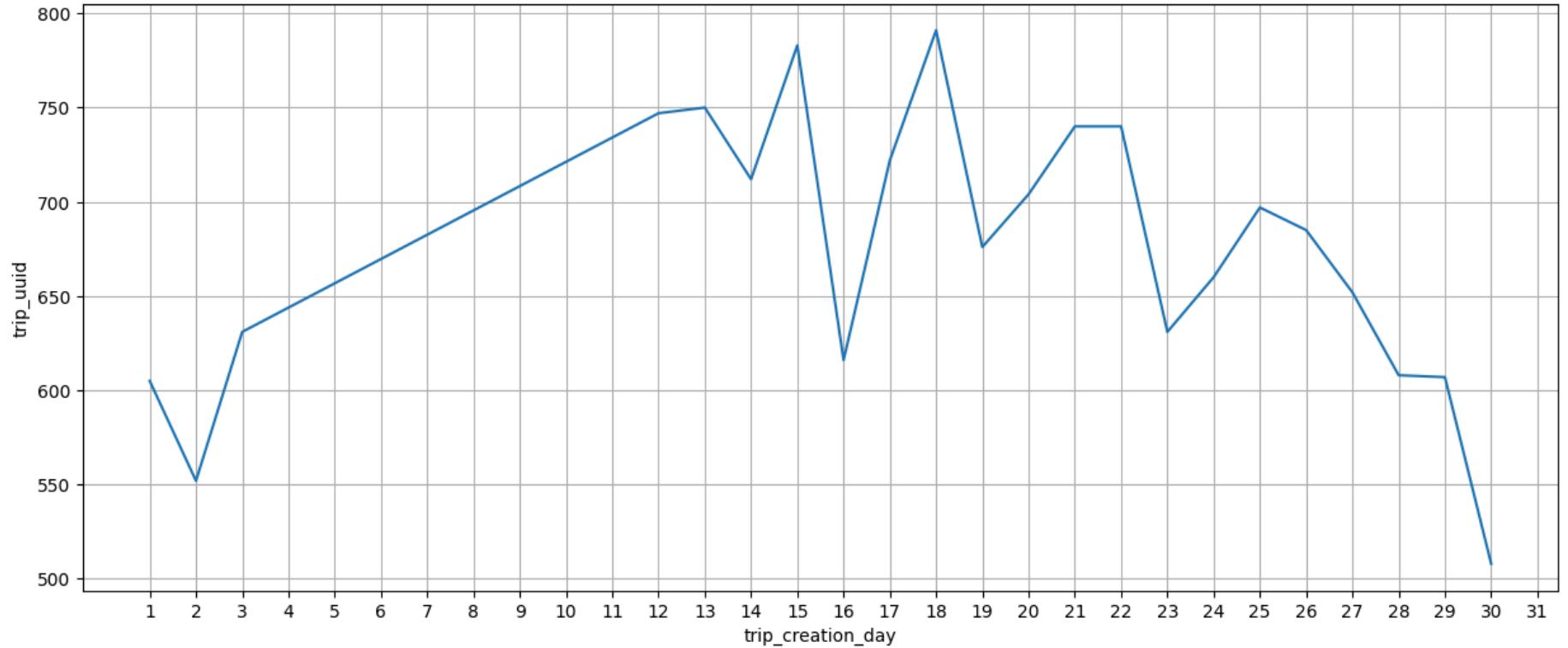
```
In [88]: df_day = data_2.groupby(by = 'trip_creation_day')['trip_uuid'].count().to_frame().reset_index()  
df_day.head()
```

```
Out[88]:    trip_creation_day  trip_uuid
```

0	1	605
1	2	552
2	3	631
3	12	747
4	13	750

```
In [89]: plt.figure(figsize = (15, 6))
sns.lineplot(data = df_day,
              x = df_day['trip_creation_day'],
              y = df_day['trip_uuid'],
              markers = 'o')
plt.xticks(np.arange(1, 32))
plt.grid('both')
plt.plot()
```

```
Out[89]: []
```



```
In [90]: #It can be inferred from the above plot that most of the trips are created in the mid of the month.  
#That means customers usually make more orders in the mid of the month.
```

```
In [91]: #I am interested to know how many trips are created for different weeks  
data_2['trip_creation_week'].unique()
```

```
Out[91]: array([37, 38, 39, 40], dtype=int8)
```

```
In [93]: df_week = data_2.groupby(by = 'trip_creation_week')['trip_uuid'].count().to_frame().reset_index()  
df_week.head()
```

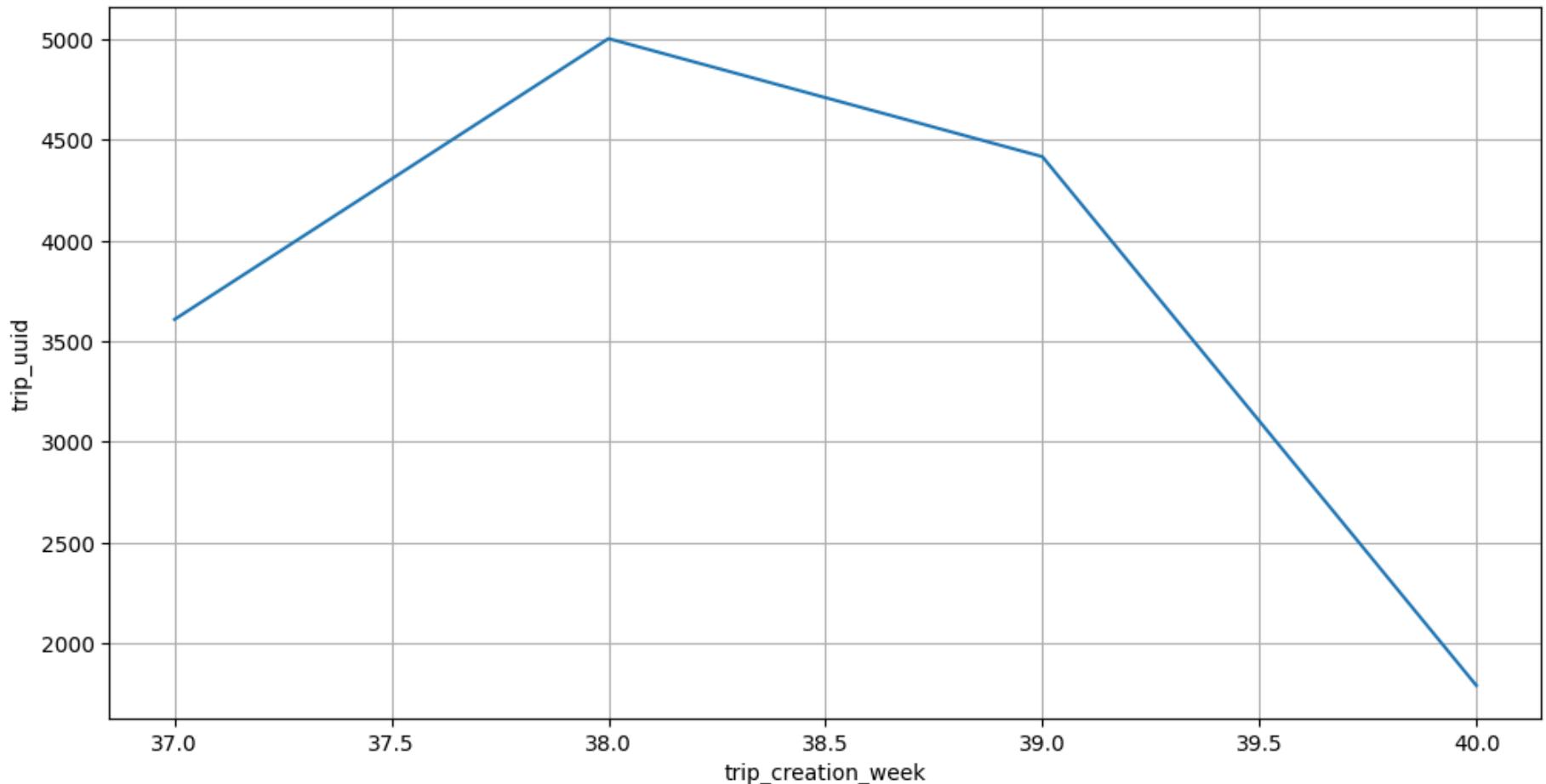
```
Out[93]:    trip_creation_week  trip_uuid
```

0	37	3608
1	38	5004
2	39	4417
3	40	1788

```
In [94]: plt.figure(figsize = (12, 6))
```

```
  sns.lineplot(data = df_week,
                x = df_week['trip_creation_week'],
                y = df_week['trip_uuid'],
                markers = 'o')
plt.grid('both')
plt.plot()
```

```
Out[94]: []
```



```
In [95]: #It can be inferred from the above plot that most of the trips are created in the 38th week.
```

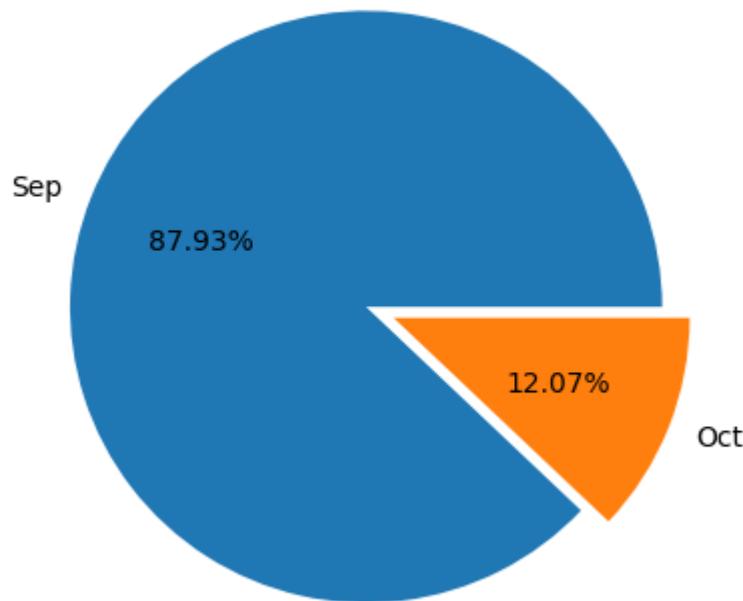
```
In [96]: #I am interested to know how many trips are created in the given two months  
df_month = data_2.groupby(by = 'trip_creation_month')['trip_uuid'].count().to_frame().reset_index()  
df_month['perc'] = np.round(df_month['trip_uuid'] * 100/ df_month['trip_uuid'].sum(), 2)  
df_month.head()
```

```
Out[96]:    trip_creation_month  trip_uuid   perc
```

	trip_creation_month	trip_uuid	perc
0	9	13029	87.93
1	10	1788	12.07

```
In [97]: plt.pie(x = df_month['trip_uuid'],
               labels = ['Sep', 'Oct'],
               explode = [0, 0.1],
               autopct = '%.2f%%')
plt.plot()
```

```
Out[97]: []
```



```
In [98]: #I am interested to know the distribution of trip data for the orders
```

```
df_data = data_2.groupby(by = 'data')['trip_uuid'].count().to_frame().reset_index()
```

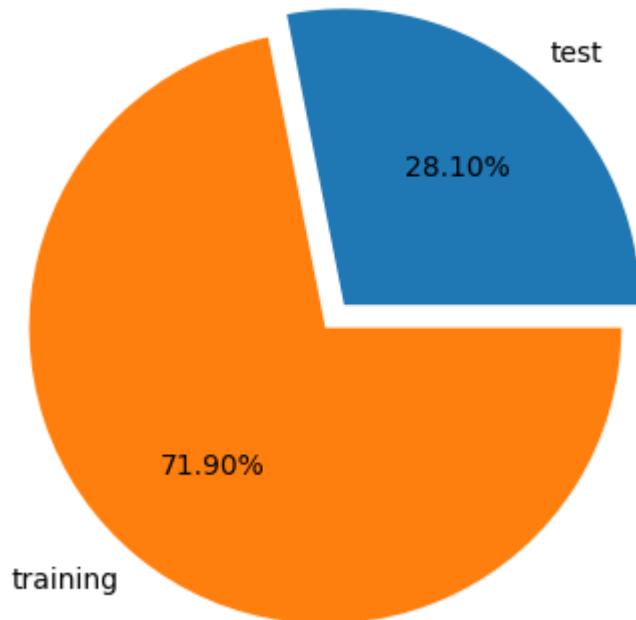
```
df_data['perc'] = np.round(df_data['trip_uuid'] * 100 / df_data['trip_uuid'].sum(), 2)
df_data.head()
```

Out[98]:

	data	trip_uuid	perc
0	test	4163	28.1
1	training	10654	71.9

```
In [99]: plt.pie(x = df_data['trip_uuid'],
               labels = df_data['data'],
               explode = [0, 0.1],
               autopct = '%.2f%')
plt.plot()
```

Out[99]: []



```
In [100...]
```

```
#I am interested to know the distribution of route types for the orders
df_route = data_2.groupby(by = 'route_type')['trip_uuid'].count().to_frame().reset_index()
df_route['perc'] = np.round(df_route['trip_uuid'] * 100/ df_route['trip_uuid'].sum(), 2)
df_route.head()
```

```
Out[100...]
```

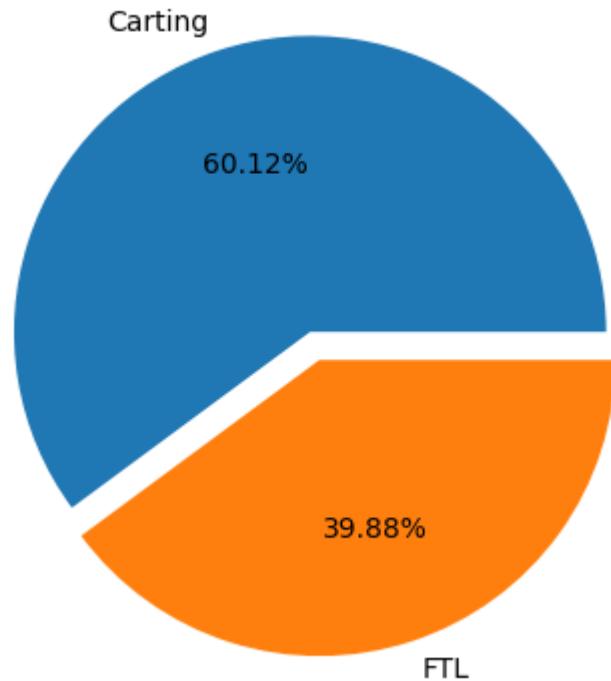
	route_type	trip_uuid	perc
0	Carting	8908	60.12
1	FTL	5909	39.88

```
In [101...]
```

```
plt.pie(x = df_route['trip_uuid'],
         labels = ['Carting', 'FTL'],
         explode = [0, 0.1],
         autopct = '%.2f%%')
plt.plot()
```

```
Out[101...]
```

```
[]
```



In [102...]

```
#I am interested to know what is the distribution of number of trips created from different states
df_source_state = data_2.groupby(by = 'source_state')['trip_uuid'].count().to_frame().reset_index()
df_source_state['perc'] = np.round(df_source_state['trip_uuid'] * 100/ df_source_state['trip_uuid'].sum(), 2)
df_source_state = df_source_state.sort_values(by = 'trip_uuid', ascending = False)
df_source_state.head()
```

Out[102...]

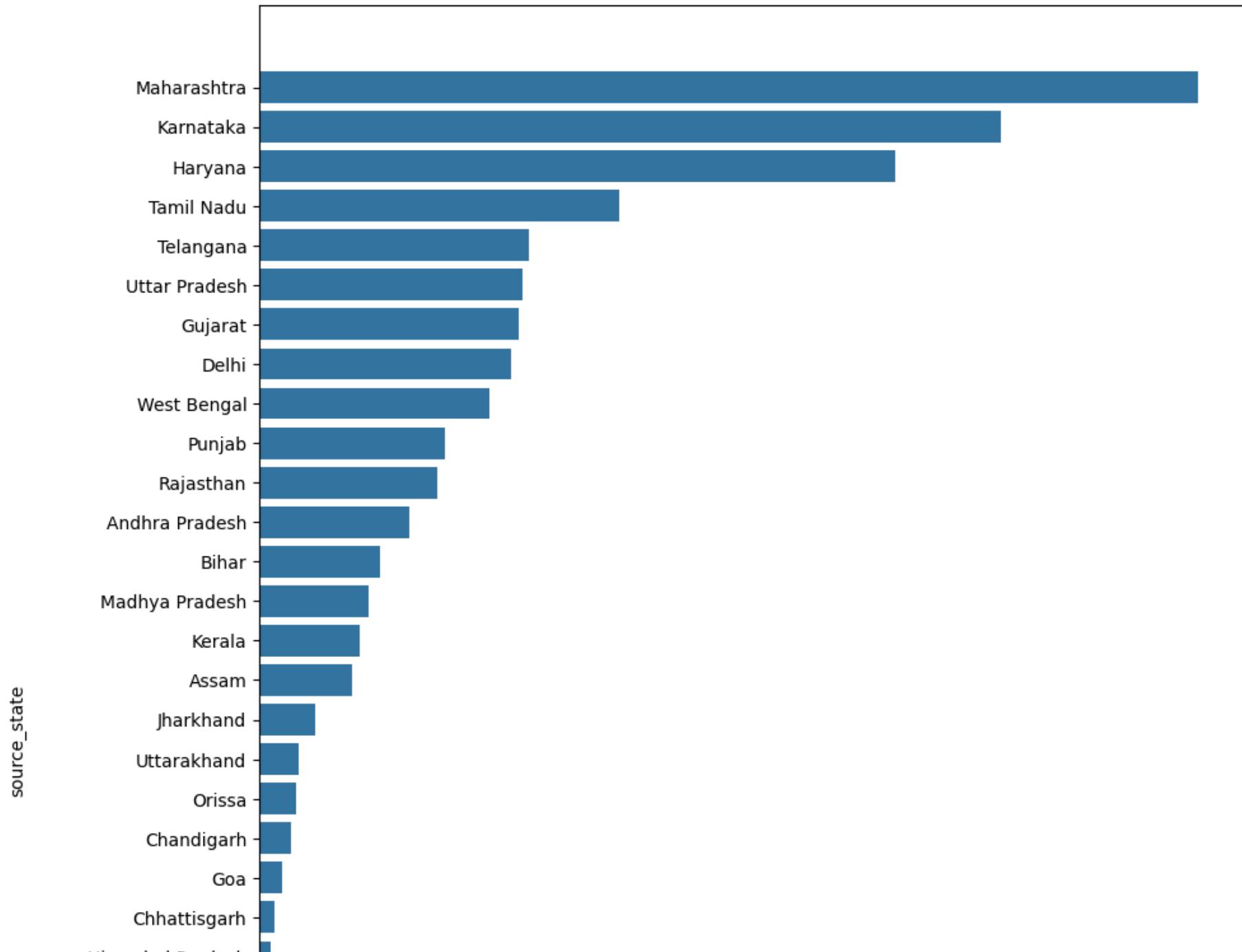
	source_state	trip_uuid	perc
17	Maharashtra	2714	18.32
14	Karnataka	2143	14.46
10	Haryana	1838	12.40
24	Tamil Nadu	1039	7.01
25	Telangana	781	5.27

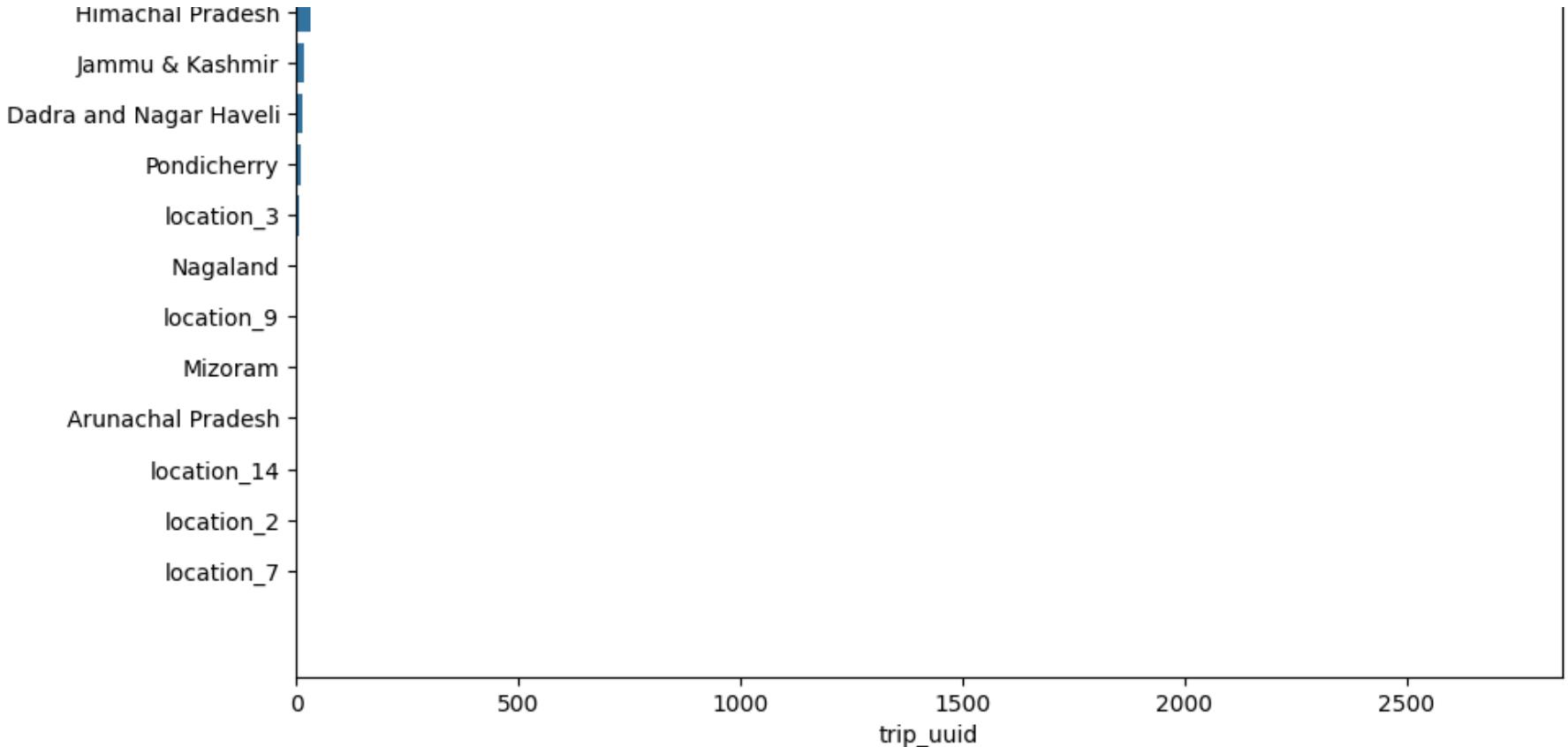
In [103...]

```
plt.figure(figsize = (10, 15))
sns.barplot(data = df_source_state,
             x = df_source_state['trip_uuid'],
             y = df_source_state['source_state'])
plt.plot()
```

Out[103...]

[]





```
In [105...]: #It can be seen in the above plot that maximum trips originated from Maharashtra state followed by Karnataka and Haryana.  
#That means that the seller base is strong in these states
```

```
In [106...]: #I am interested to know top 30 cities based on the number of trips created from different cities  
df_source_city = data_2.groupby(by = 'source_city')['trip_uuid'].count().to_frame().reset_index()  
df_source_city['perc'] = np.round(df_source_city['trip_uuid'] * 100/ df_source_city['trip_uuid'].sum(), 2)  
df_source_city = df_source_city.sort_values(by = 'trip_uuid', ascending = False)[:30]  
df_source_city
```

Out[106...]

	source_city	trip_uuid	perc
439	Mumbai	1442	9.73
237	Gurgaon	1165	7.86
169	Delhi	883	5.96
79	Bengaluru	726	4.90
100	Bhiwandi	697	4.70
58	Bangalore	648	4.37
136	Chennai	568	3.83
264	Hyderabad	524	3.54
516	Pune	480	3.24
357	Kolkata	356	2.40
610	Sonipat	276	1.86
2	Ahmedabad	274	1.85
133	Chandigarh	273	1.84
270	Jaipur	259	1.75
201	Faridabad	227	1.53
447	Muzaffarpur	159	1.07
382	Ludhiana	158	1.07
320	Kanpur	145	0.98
621	Surat	140	0.94
473	Noida	129	0.87
102	Bhopal	125	0.84
240	Guwahati	118	0.80

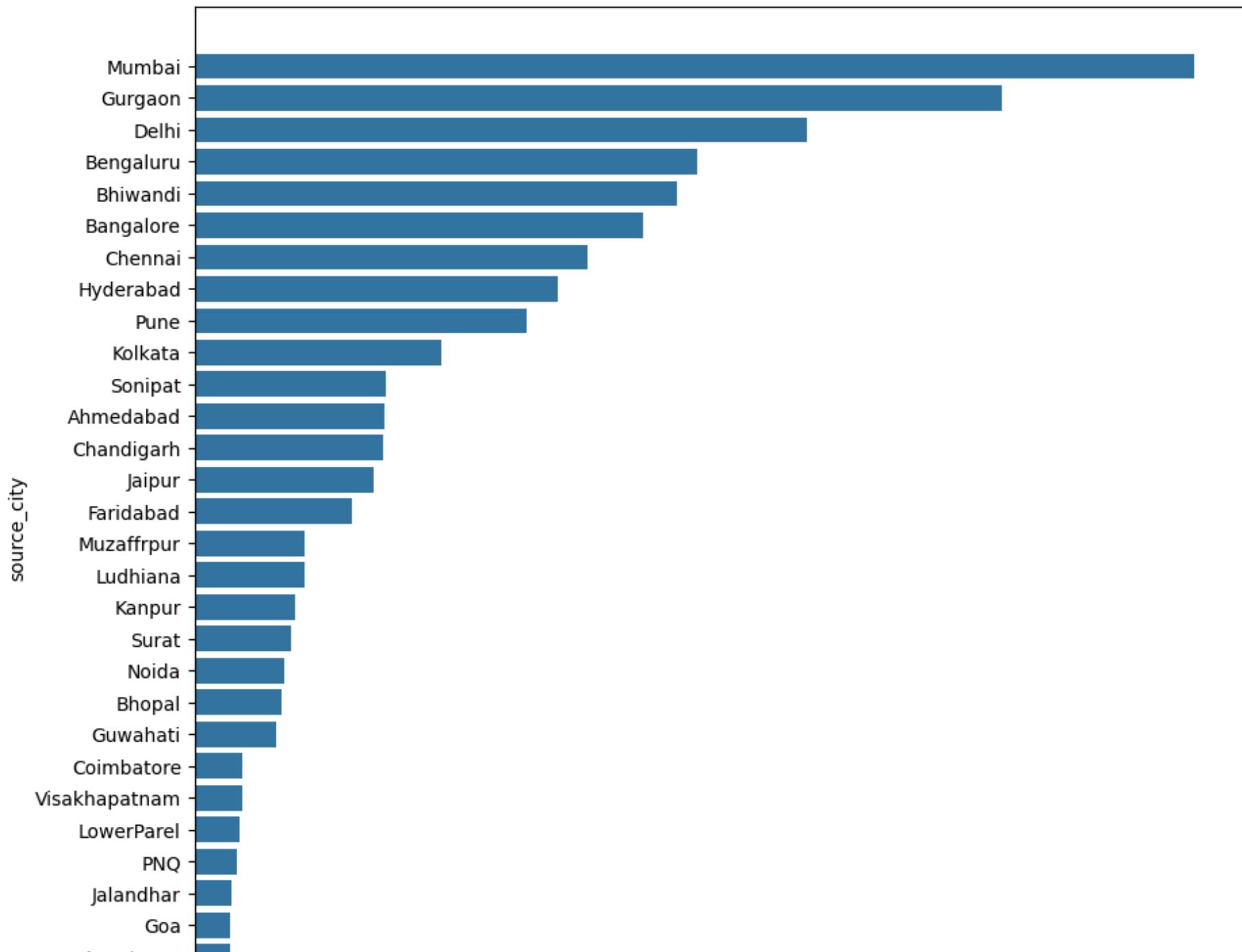
	source_city	trip_uuid	perc
154	Coimbatore	69	0.47
679	Visakhapatnam	69	0.47
380	LowerParel	65	0.44
477	PNQ	62	0.42
273	Jalandhar	54	0.36
220	Goa	52	0.35
25	Anantapur	51	0.34
261	Hubli	47	0.32

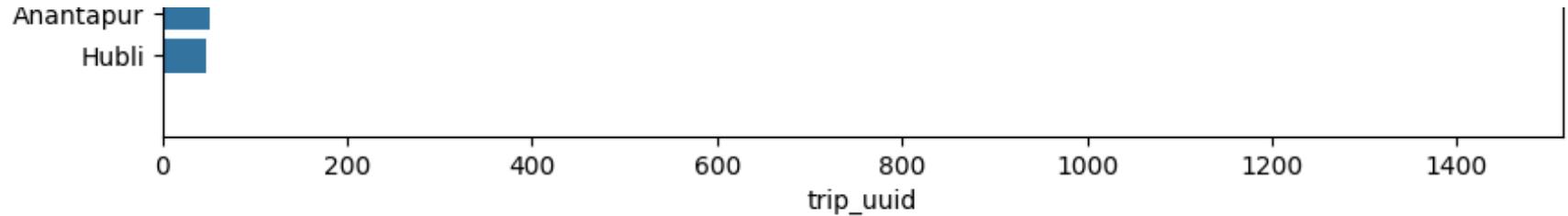
In [107...]

```
plt.figure(figsize = (10, 10))
sns.barplot(data = df_source_city,
             x = df_source_city['trip_uuid'],
             y = df_source_city['source_city'])
plt.plot()
```

Out[107...]

[]





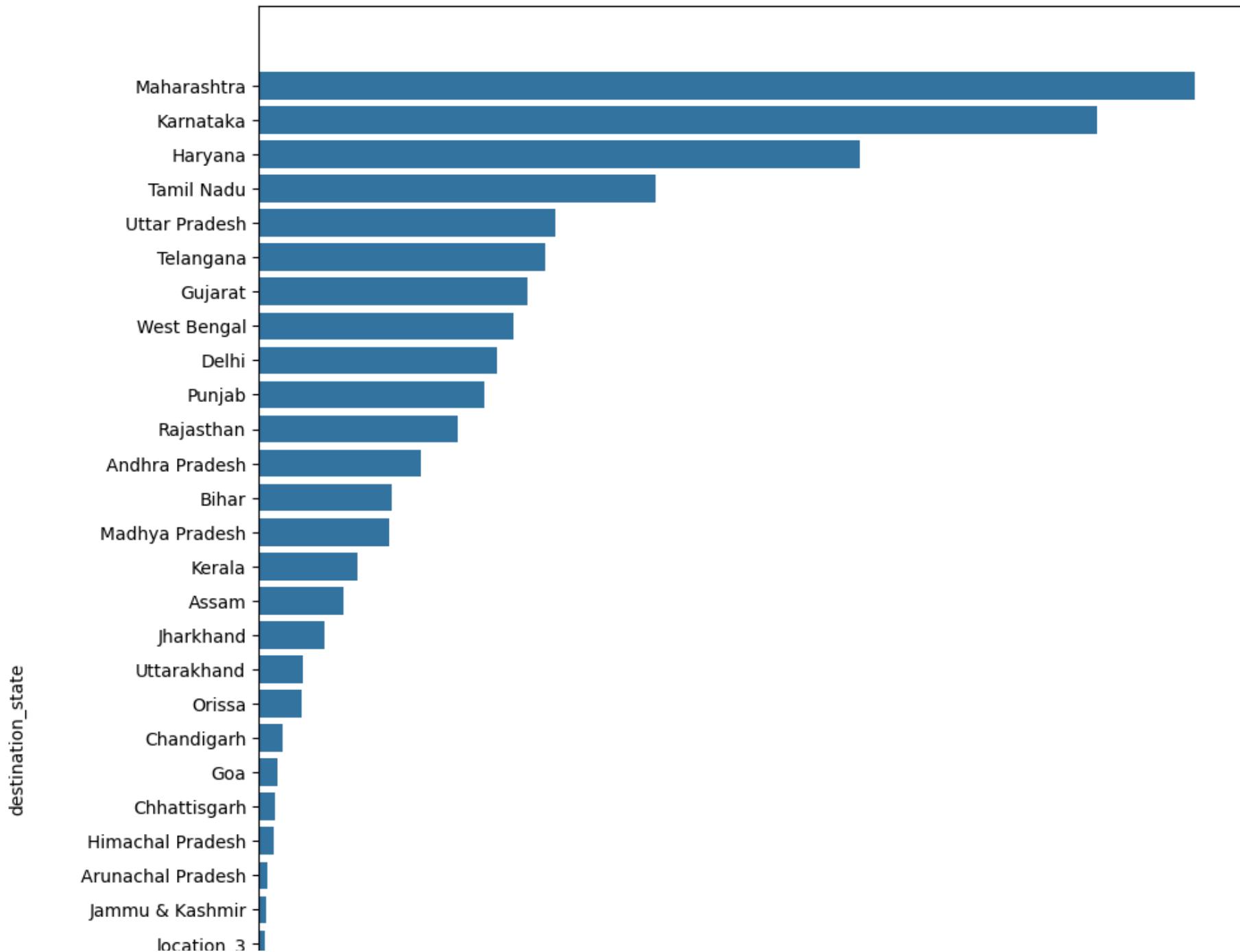
```
In [109...]: #It can be seen in the above plot that maximum trips originated from Mumbai city followed by Gurgaon Delhi, Bengaluru and Bhiwani  
#That means that the seller base is strong in these cities.
```

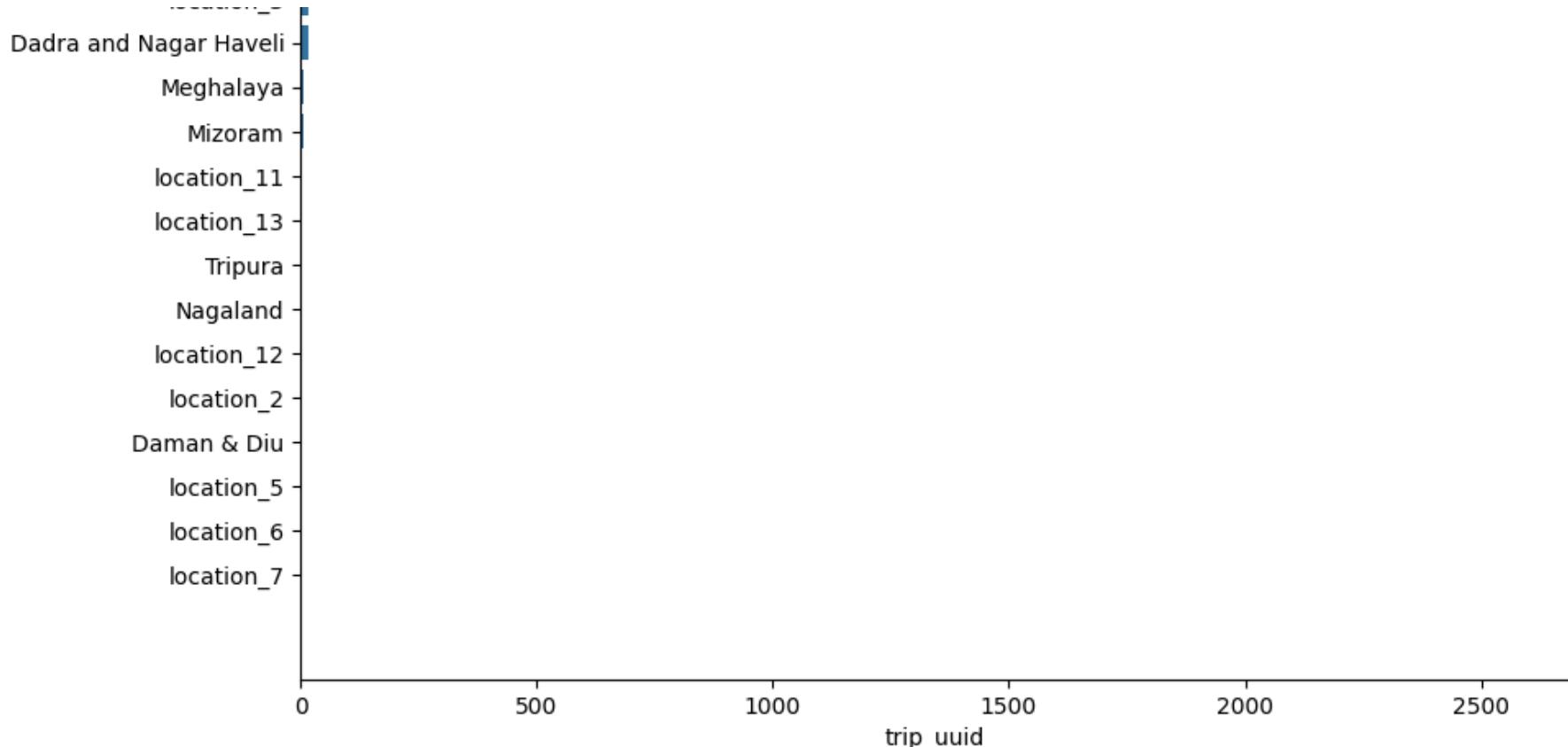
```
In [110...]: df_destination_state = data_2.groupby(by = 'destination_state')['trip_uuid'].count().to_frame().reset_index()  
df_destination_state['perc'] = np.round(df_destination_state['trip_uuid'] * 100/ df_destination_state['trip_uuid'].sum(), 2)  
df_destination_state = df_destination_state.sort_values(by = 'trip_uuid', ascending = False)  
df_destination_state.head()
```

```
Out[110...]: destination_state  trip_uuid  perc  
18 Maharashtra  2561  17.28  
15 Karnataka  2294  15.48  
11 Haryana  1643  11.09  
25 Tamil Nadu  1084  7.32  
28 Uttar Pradesh  811  5.47
```

```
In [111...]: plt.figure(figsize = (10, 15))  
sns.barplot(data = df_destination_state,  
            x = df_destination_state['trip_uuid'],  
            y = df_destination_state['destination_state'])  
plt.plot()
```

```
Out[111...]: []
```





```
In [112]: #It can be seen in the above plot that maximum trips ended in Maharashtra state followed by Karnataka, Haryana, Tamil Nadu and  
#That means that the number of orders placed in these states is significantly high in these states.
```

```
In [114]: #I am interested to know top 30 cities based on the number of trips ended in different cities  
df_destination_city = data_2.groupby(by = 'destination_city')['trip_uuid'].count().to_frame().reset_index()  
df_destination_city['perc'] = np.round(df_destination_city['trip_uuid'] * 100/ df_destination_city['trip_uuid'].sum(), 2)  
df_destination_city = df_destination_city.sort_values(by = 'trip_uuid', ascending = False)[:30]  
df_destination_city
```

Out[114...]

	destination_city	trip_uuid	perc
515	Mumbai	1548	10.45
96	Bengaluru	975	6.58
282	Gurgaon	936	6.32
200	Delhi	778	5.25
163	Chennai	595	4.02
72	Bangalore	551	3.72
308	Hyderabad	503	3.39
115	Bhiwandi	434	2.93
418	Kolkata	384	2.59
158	Chandigarh	339	2.29
724	Sonipat	322	2.17
612	Pune	317	2.14
4	Ahmedabad	265	1.79
242	Faridabad	244	1.65
318	Jaipur	205	1.38
371	Kanpur	148	1.00
117	Bhopal	139	0.94
559	PNQ	122	0.82
739	Surat	117	0.79
552	Noida	106	0.72
521	Muzaffrpur	102	0.69
284	Guwahati	98	0.66

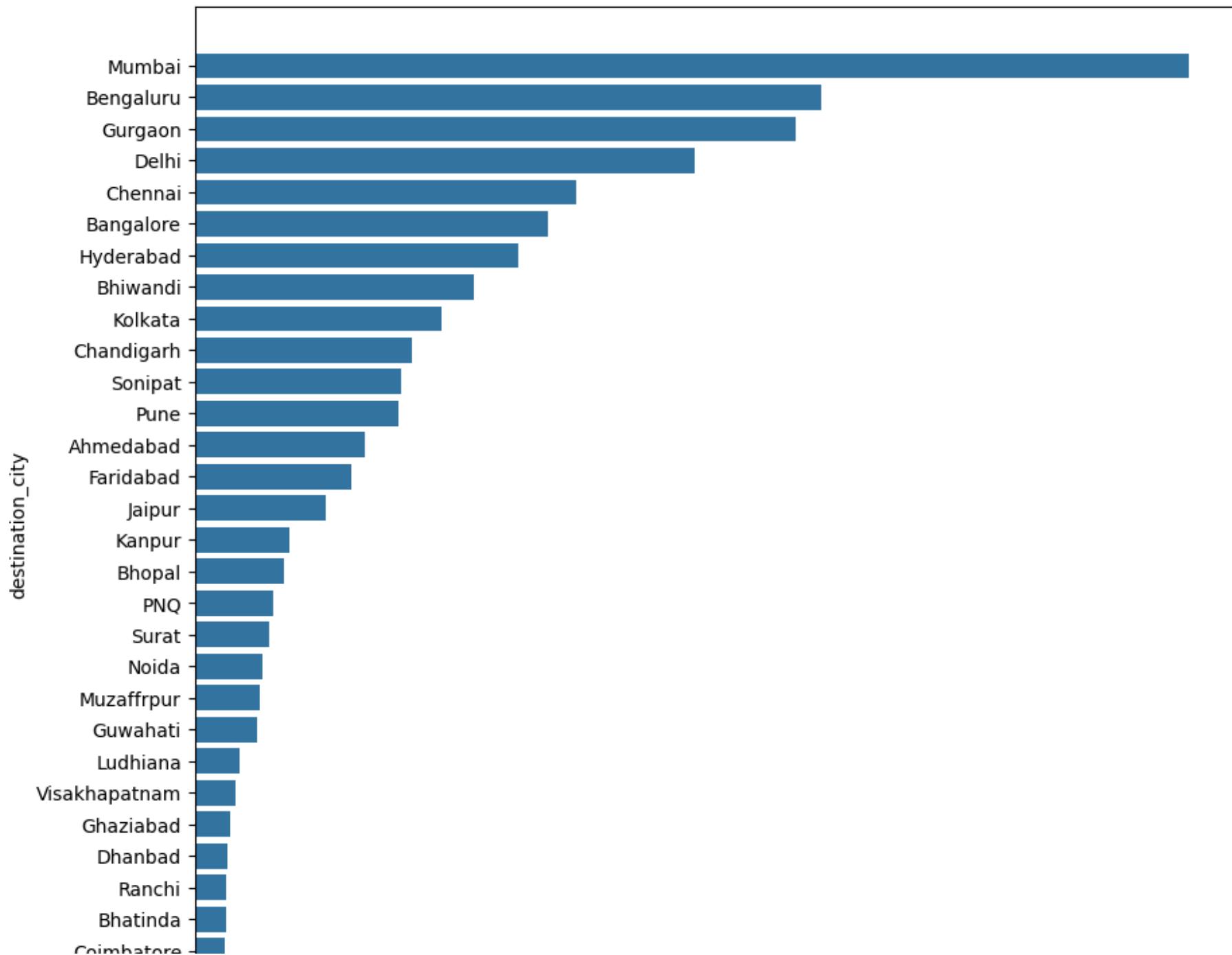
	destination_city	trip_uuid	perc
448	Ludhiana	70	0.47
797	Visakhapatnam	64	0.43
259	Ghaziabad	56	0.38
208	Dhanbad	50	0.34
639	Ranchi	49	0.33
110	Bhatinda	48	0.32
183	Coimbatore	47	0.32
9	Akola	45	0.30

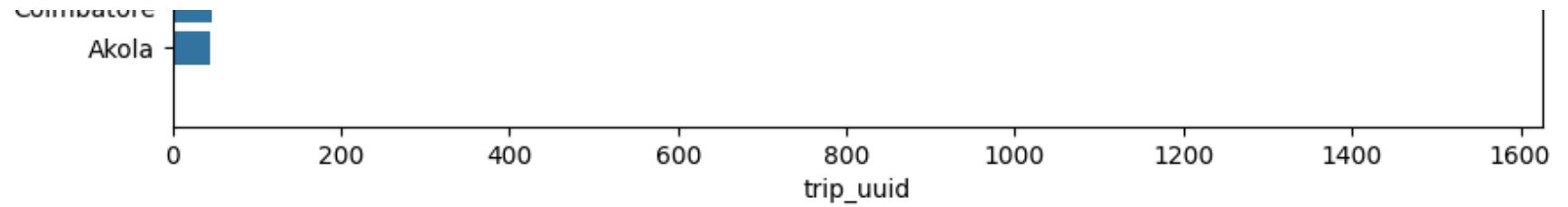
In [115...]

```
plt.figure(figsize = (10, 10))
sns.barplot(data = df_destination_city,
             x = df_destination_city['trip_uuid'],
             y = df_destination_city['destination_city'])
plt.plot()
```

Out[115...]

[]

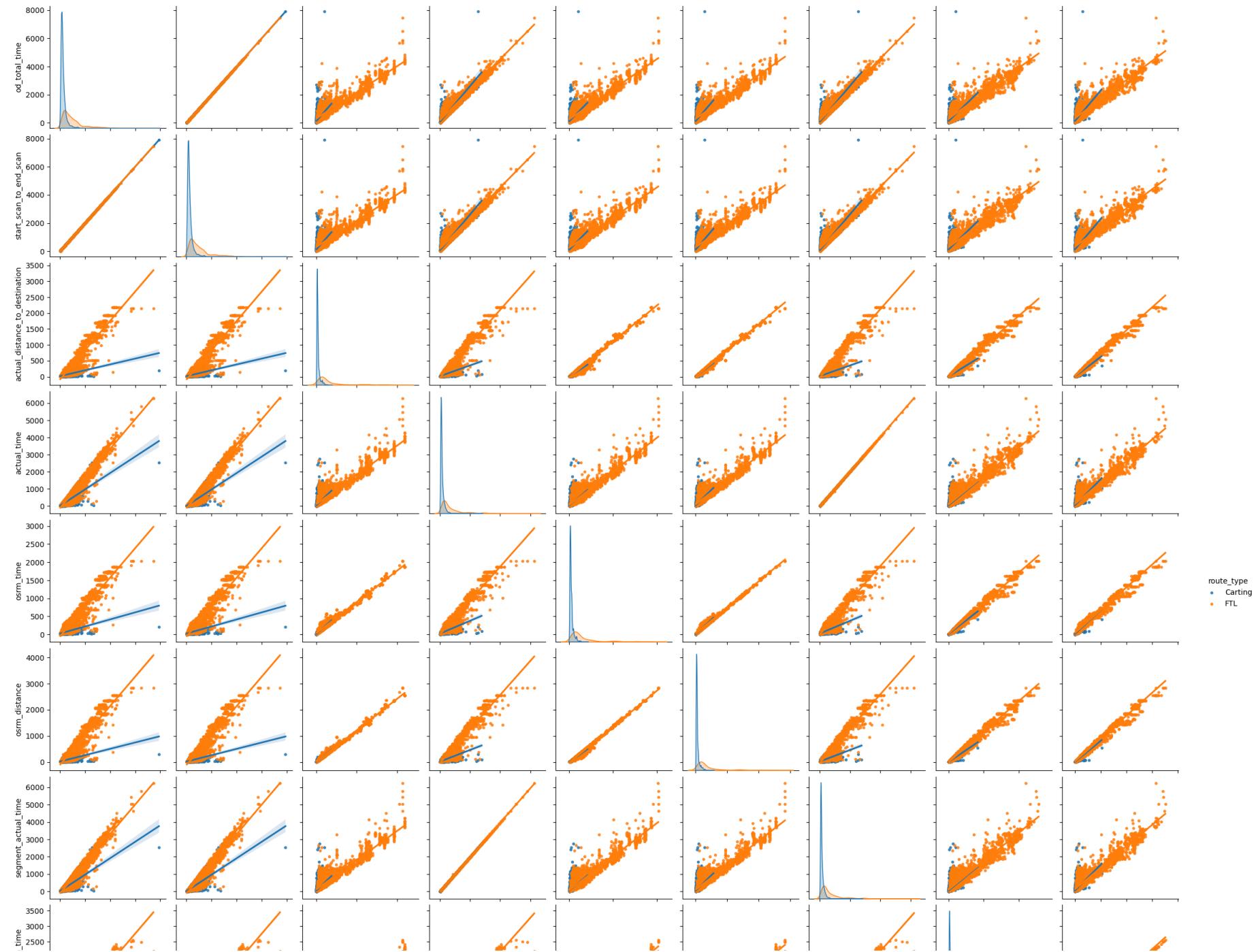


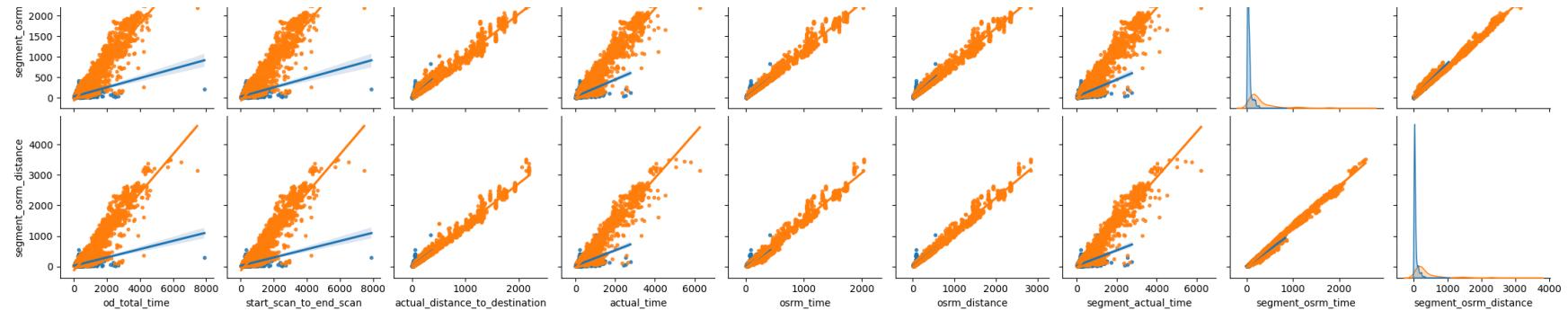


```
In [116...]: #It can be seen in the above plot that maximum trips ended in Mumbai city followed by Bengaluru, Gurgaon, Delhi and Chennai.  
#That means that the number of orders placed in these cities is significantly high.
```

```
In [118...]: numerical_columns = ['od_total_time', 'start_scan_to_end_scan', 'actual_distance_to_destination',  
                      'actual_time', 'osrm_time', 'osrm_distance', 'segment_actual_time',  
                      'segment_osrm_time', 'segment_osrm_distance']  
sns.pairplot(data = data_2,  
              vars = numerical_columns,  
              kind = 'reg',  
              hue = 'route_type',  
              markers = '.')  
plt.plot()
```

```
Out[118...]: []
```





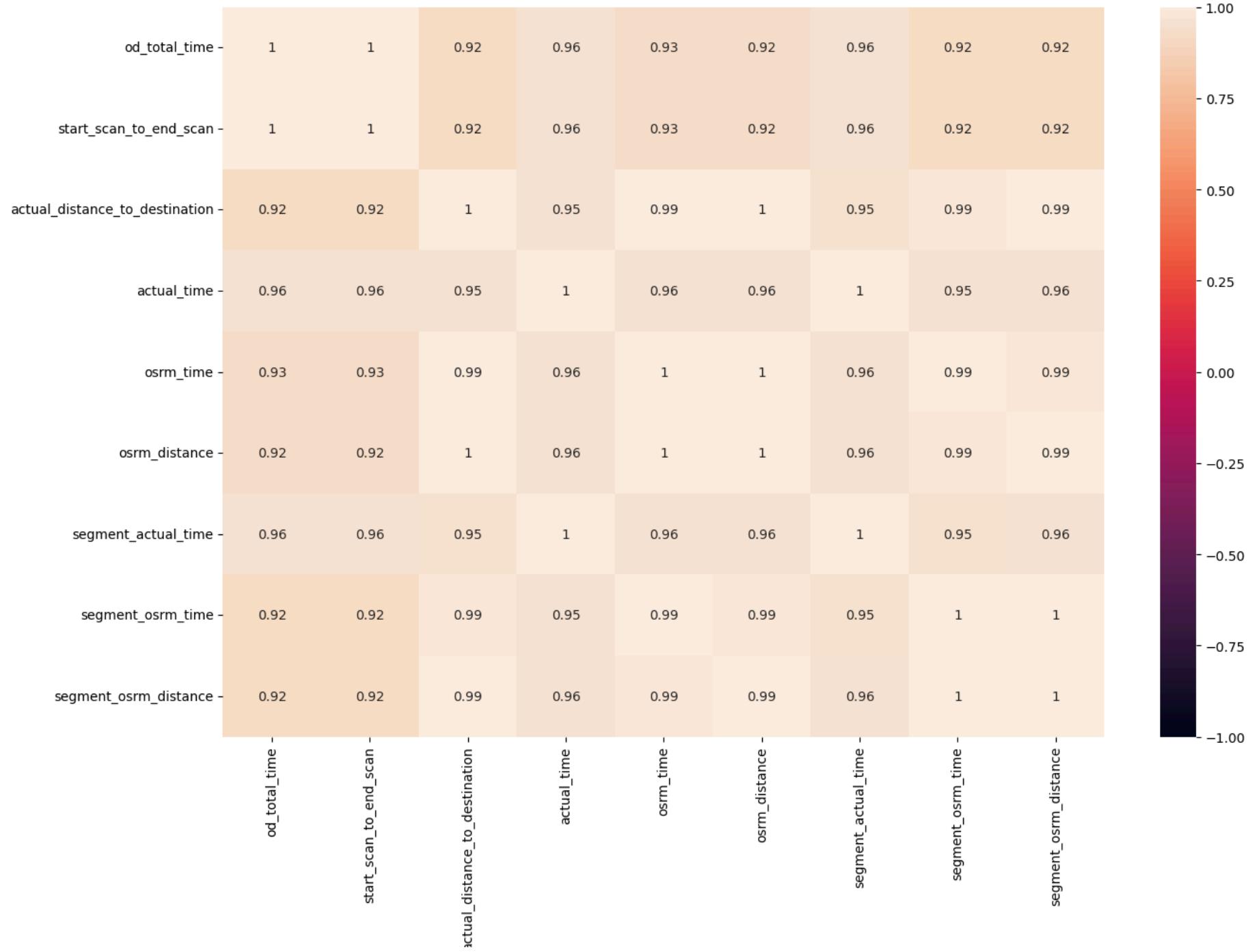
In [119...]
`df_corr = data_2[numerical_columns].corr()
df_corr`

Out[119...]

	od_total_time	start_scan_to_end_scan	actual_distance_to_destination	actual_time	osrm_time	osrm_distance	
od_total_time	1.000000	0.999999		0.918222	0.961094	0.926516	0.924219
start_scan_to_end_scan	0.999999	1.000000		0.918308	0.961147	0.926571	0.924299
actual_distance_to_destination	0.918222	0.918308	1.000000	0.953757	0.993561	0.997264	
actual_time	0.961094	0.961147	0.953757	1.000000	0.958593	0.959214	
osrm_time	0.926516	0.926571	0.993561	0.958593	1.000000	0.997580	
osrm_distance	0.924219	0.924299	0.997264	0.959214	0.997580	1.000000	
segment_actual_time	0.961119	0.961171	0.952821	0.999989	0.957765	0.958353	
segment_osrm_time	0.918490	0.918561	0.987538	0.953872	0.993259	0.991798	
segment_osrm_distance	0.919199	0.919291	0.993061	0.956967	0.991608	0.994710	

In [120...]
`plt.figure(figsize = (15, 10))
sns.heatmap(data = df_corr, vmin = -1, vmax = 1, annot = True)
plt.plot()`

Out[120...]
[]



In [121...]

```
#Very High Correlation (> 0.9) exists between columns all the numerical columns specified above
```

In-depth analysis and feature engineering:

Compare the difference between od_total_time and start_scan_to_end_scan. Do hypothesis testing/ Visual analysis to check

STEP-1 : Set up Null Hypothesis

Null Hypothesis (H0) - od_total_time (Total Trip Time) and start_scan_to_end_scan (Expected total trip time) are same. Alternate Hypothesis (HA) - od_total_time (Total Trip Time) and start_scan_to_end_scan (Expected total trip time) are different.

STEP-2 : Checking for basic assumptions for the hypothesis

Distribution check using QQ Plot Homogeneity of Variances using Lavene's test

STEP-3: Define Test statistics; Distribution of T under H0.

If the assumptions of T Test are met then we can proceed performing T Test for independent samples else we will perform the non parametric test equivalent to T Test for independent sample i.e., Mann-Whitney U rank test for two independent samples.

STEP-4: Compute the p-value and fix value of alpha.

We set our alpha to be 0.05

STEP-5: Compare p-value and alpha.

Based on p-value, we will accept or reject H0.

p-val > alpha : Accept H0

p-val < alpha : Reject H0

```
In [125...]: data_2[['od_total_time', 'start_scan_to_end_scan']].describe()
```

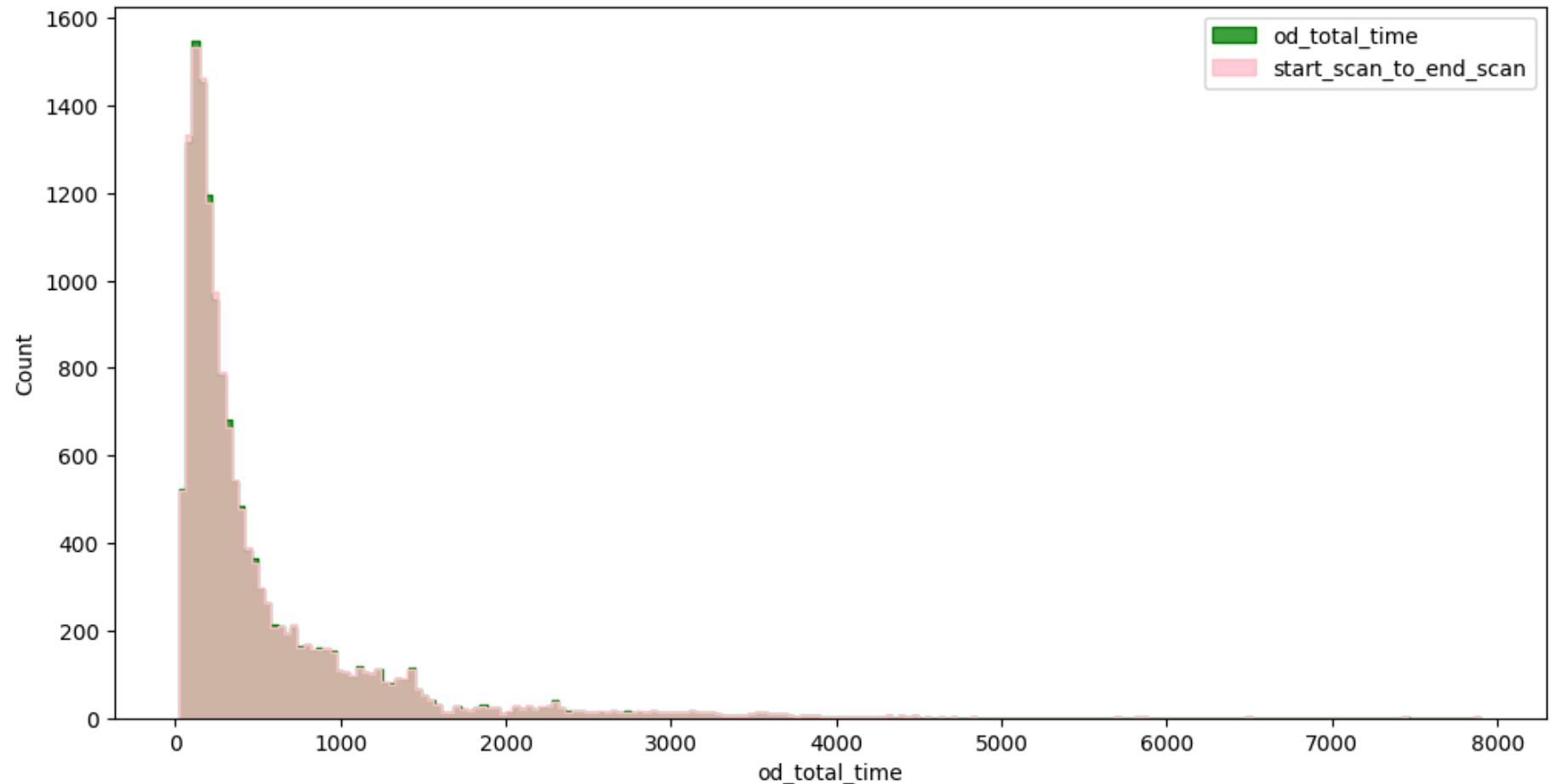
```
Out[125...]:
```

	od_total_time	start_scan_to_end_scan
count	14817.000000	14817.000000
mean	531.697630	530.810016
std	658.868223	658.705957
min	23.460000	23.000000
25%	149.930000	149.000000
50%	280.770000	280.000000
75%	638.200000	637.000000
max	7898.550000	7898.000000

```
In [127...]: #Visual Tests to know if the samples follow normal distribution
```

```
plt.figure(figsize = (12, 6))
sns.histplot(data_2['od_total_time'], element = 'step', color = 'green')
sns.histplot(data_2['start_scan_to_end_scan'], element = 'step', color = 'pink')
plt.legend(['od_total_time', 'start_scan_to_end_scan'])
plt.plot()
```

```
Out[127...]: []
```

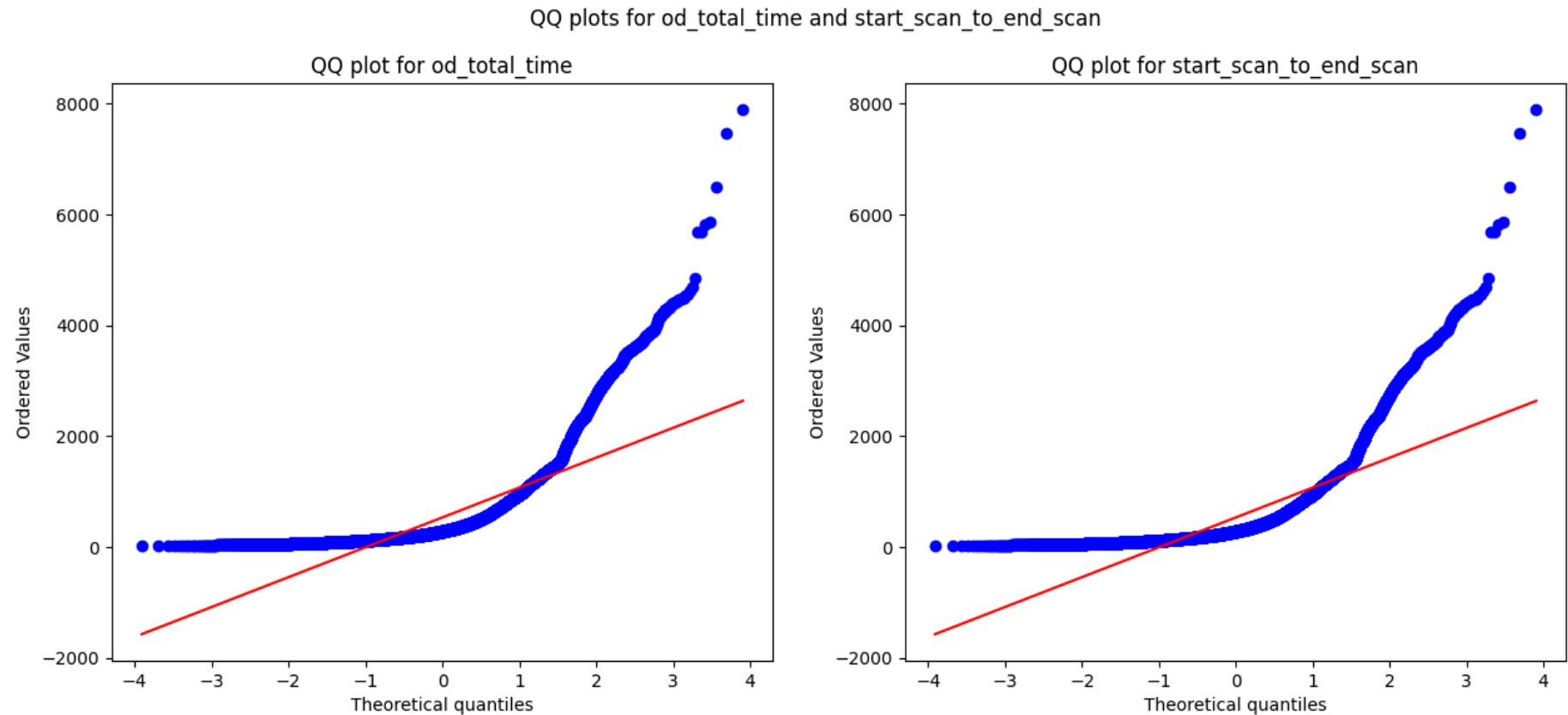


```
In [129]: import scipy.stats as spy
```

```
In [130]: #Distribution check using QQ Plot
plt.figure(figsize = (15, 6))
plt.subplot(1, 2, 1)
plt.suptitle('QQ plots for od_total_time and start_scan_to_end_scan')
spy.probplot(data_2['od_total_time'], plot = plt, dist = 'norm')
plt.title('QQ plot for od_total_time')
plt.subplot(1, 2, 2)
spy.probplot(data_2['start_scan_to_end_scan'], plot = plt, dist = 'norm')
```

```
plt.title('QQ plot for start_scan_to_end_scan')  
plt.plot()
```

Out[130...]: []



It can be seen from the above plots that the samples do not come from normal distribution.

Applying Shapiro-Wilk test for normality

H0 : The sample follows normal distribution

H1 : The sample does not follow normal distribution

alpha = 0.05

```
In [131...]: #Test Statistics : Shapiro-Wilk test for normality

test_stat, p_value = spy.shapiro(data_2['od_total_time'].sample(5000))
print('p-value', p_value)
if p_value < 0.05:
    print('The sample does not follow normal distribution')
else:
    print('The sample follows normal distribution')

p-value 0.0
The sample does not follow normal distribution

In [132...]: test_stat, p_value = spy.shapiro(data_2['start_scan_to_end_scan'].sample(5000))
print('p-value', p_value)
if p_value < 0.05:
    print('The sample does not follow normal distribution')
else:
    print('The sample follows normal distribution')

p-value 0.0
The sample does not follow normal distribution

In [133...]: #Transforming the data using boxcox transformation to check if the transformed data follows normal distribution.
transformed_od_total_time = spy.boxcox(data_2['od_total_time'])[0]
test_stat, p_value = spy.shapiro(transformed_od_total_time)
print('p-value', p_value)
if p_value < 0.05:
    print('The sample does not follow normal distribution')
else:
    print('The sample follows normal distribution')

p-value 7.172770042757021e-25
The sample does not follow normal distribution

In [134...]: transformed_start_scan_to_end_scan = spy.boxcox(data_2['start_scan_to_end_scan'])[0]
test_stat, p_value = spy.shapiro(transformed_start_scan_to_end_scan)
print('p-value', p_value)
if p_value < 0.05:
    print('The sample does not follow normal distribution')
else:
    print('The sample follows normal distribution')
```

```
p-value 1.0471322892609475e-24  
The sample does not follow normal distribution
```

```
In [135...]: #1. Even after applying the boxcox transformation on each of the "od_total_time" and "start_scan_to_end_scan" columns, the dis  
#2. Homogeneity of Variances using Lavene's test
```

```
In [136...]: # Null Hypothesis(H0) - Homogenous Variance  
# Alternate Hypothesis(HA) - Non Homogenous Variance  
  
test_stat, p_value = spy.levene(data_2['od_total_time'], data_2['start_scan_to_end_scan'])  
print('p-value', p_value)  
if p_value < 0.05:  
    print('The samples do not have Homogenous Variance')  
else:  
    print('The samples have Homogenous Variance ')
```

```
p-value 0.9668007217581142  
The samples have Homogenous Variance
```

```
In [137...]: #Since the samples are not normally distributed, T-Test cannot be applied here, we can perform its non parametric equivalent t  
test_stat, p_value = spy.mannwhitneyu(data_2['od_total_time'], data_2['start_scan_to_end_scan'])  
print('P-value :',p_value)
```

```
P-value : 0.7815123224221716
```

```
In [138...]: # Since p-value > alpha therefore it can be concluded that od_total_time and start_scan_to_end_scan are similar
```

Do hypothesis testing / visual analysis between actual_time aggregated value and OSRM time aggregated value (aggregated values are the values you'll get after merging the rows on the basis of trip_uuid)

```
In [139...]: data_2[['actual_time', 'osrm_time']].describe()
```

Out[139...]

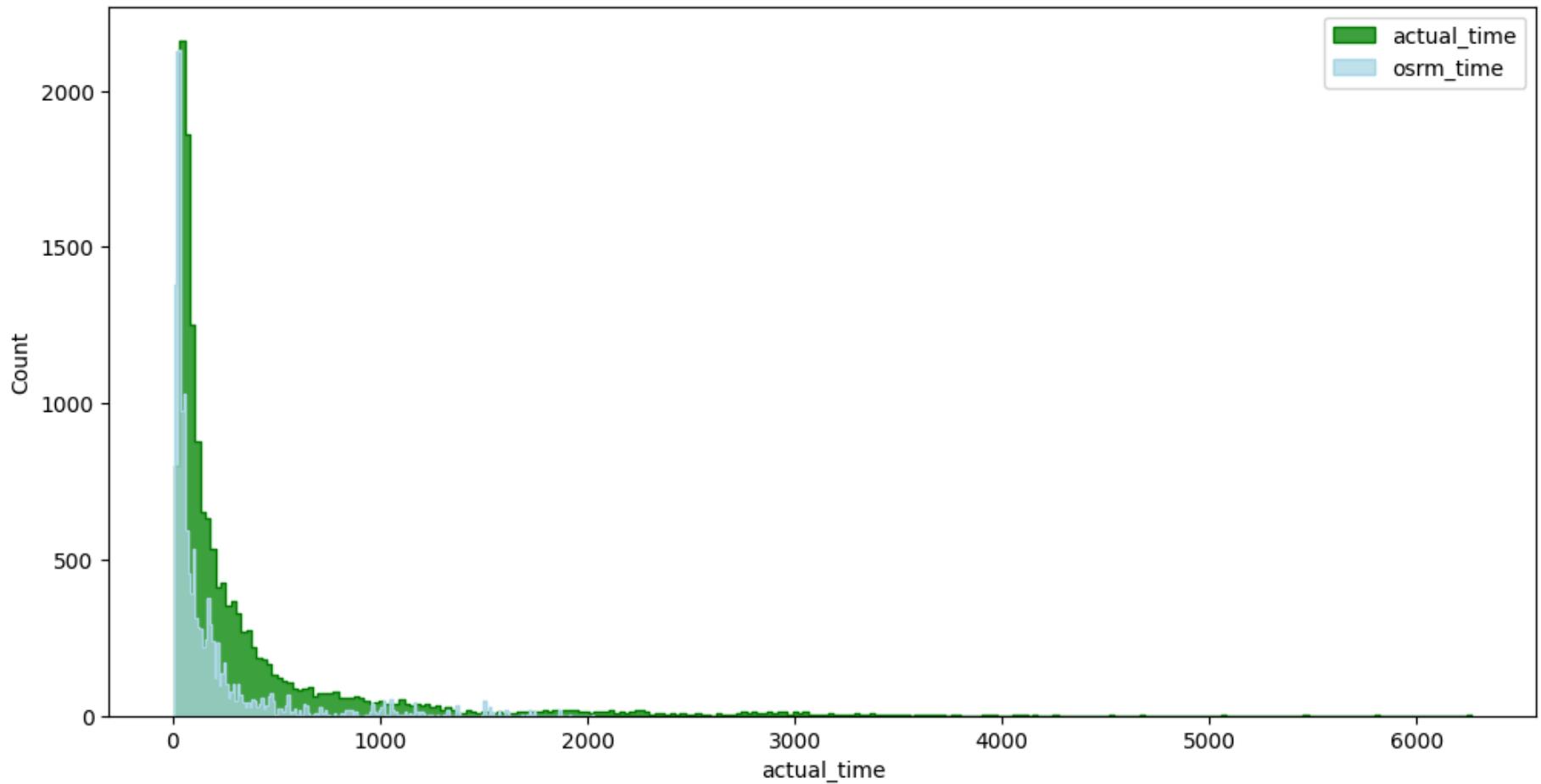
	actual_time	osrm_time
count	14817.000000	14817.000000
mean	357.143768	161.384018
std	561.396118	271.360992
min	9.000000	6.000000
25%	67.000000	29.000000
50%	149.000000	60.000000
75%	370.000000	168.000000
max	6265.000000	2032.000000

In [140...]

```
#Visual Tests to know if the samples follow normal distribution
plt.figure(figsize = (12, 6))
sns.histplot(data_2['actual_time'], element = 'step', color = 'green')
sns.histplot(data_2['osrm_time'], element = 'step', color = 'lightblue')
plt.legend(['actual_time', 'osrm_time'])
plt.plot()
```

Out[140...]

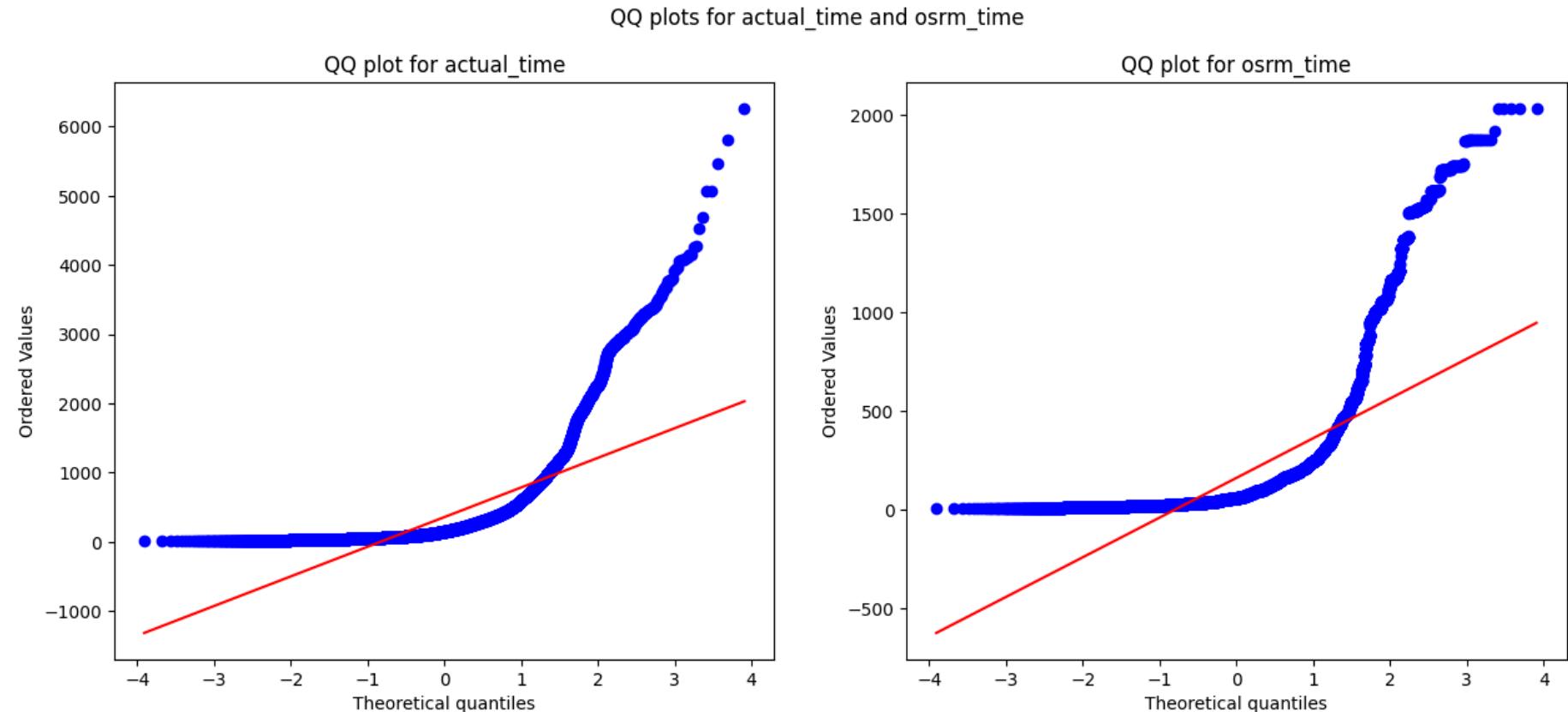
[]



In [141...]

```
# Distribution check using QQ Plot
plt.figure(figsize = (15, 6))
plt.subplot(1, 2, 1)
plt.suptitle('QQ plots for actual_time and osrm_time')
spy.probplot(data_2['actual_time'], plot = plt, dist = 'norm')
plt.title('QQ plot for actual_time')
plt.subplot(1, 2, 2)
spy.probplot(data_2['osrm_time'], plot = plt, dist = 'norm')
plt.title('QQ plot for osrm_time')
plt.plot()
```

```
Out[141...]: []
```



```
In [142...]: # It can be seen from the above plots that the samples do not come from normal distribution
```

Applying Shapiro-Wilk test for normality

H0 : The sample follows normal distribution

H1 : The sample does not follow normal distribution

alpha = 0.05

Test Statistics : Shapiro-Wilk test for normality

```
In [143...]  
test_stat, p_value = spy.shapiro(data_2['actual_time'].sample(5000))  
print('p-value', p_value)  
if p_value < 0.05:  
    print('The sample does not follow normal distribution')  
else:  
    print('The sample follows normal distribution')
```

p-value 0.0
The sample does not follow normal distribution

```
In [144...]  
test_stat, p_value = spy.shapiro(data_2['osrm_time'].sample(5000))  
print('p-value', p_value)  
if p_value < 0.05:  
    print('The sample does not follow normal distribution')  
else:  
    print('The sample follows normal distribution')
```

p-value 0.0
The sample does not follow normal distribution

```
In [145...]  
#Transforming the data using boxcox transformation to check if the transformed data follows normal distribution.  
transformed_actual_time = spy.boxcox(data_2['actual_time'])[0]  
test_stat, p_value = spy.shapiro(transformed_actual_time)  
print('p-value', p_value)  
if p_value < 0.05:  
    print('The sample does not follow normal distribution')  
else:  
    print('The sample follows normal distribution')
```

p-value 1.020620453603145e-28
The sample does not follow normal distribution

```
In [146...]  
transformed_osrm_time = spy.boxcox(data_2['osrm_time'])[0]  
test_stat, p_value = spy.shapiro(transformed_osrm_time)  
print('p-value', p_value)  
if p_value < 0.05:  
    print('The sample does not follow normal distribution')  
else:  
    print('The sample follows normal distribution')
```

```
p-value 3.5882550510138333e-35  
The sample does not follow normal distribution
```

```
In [147...]: #1. Even after applying the boxcox transformation on each of the "actual_time" and "osrm_time" columns, the distributions do not follow normal distribution  
#2. Homogeneity of Variances using Lavene's test
```

```
In [148...]:  
# Null Hypothesis(H0) - Homogenous Variance  
  
# Alternate Hypothesis(HA) - Non Homogenous Variance  
  
test_stat, p_value = spy.levene(data_2['actual_time'], data_2['osrm_time'])  
print('p-value', p_value)  
if p_value < 0.05:  
    print('The samples do not have Homogenous Variance')  
else:  
    print('The samples have Homogenous Variance ')
```

```
p-value 1.871098057987424e-220  
The samples do not have Homogenous Variance
```

```
In [149...]: #Since the samples do not follow any of the assumptions T-Test cannot be applied here, we can perform its non parametric equivalent  
test_stat, p_value = spy.mannwhitneyu(data_2['actual_time'], data_2['osrm_time'])  
print('p-value', p_value)  
if p_value < 0.05:  
    print('The samples are not similar')  
else:  
    print('The samples are similar ')
```

```
p-value 0.0  
The samples are not similar
```

```
In [150...]: #Since p-value < alpha therefore it can be concluded that actual_time and osrm_time are not similar.
```

Do hypothesis testing/ visual analysis between actual_time aggregated value and segment actual time aggregated value (aggregated values are the values you'll get after merging the rows on the basis of trip_uuid)

In [152... `data_2[['actual_time', 'segment_actual_time']].describe()`

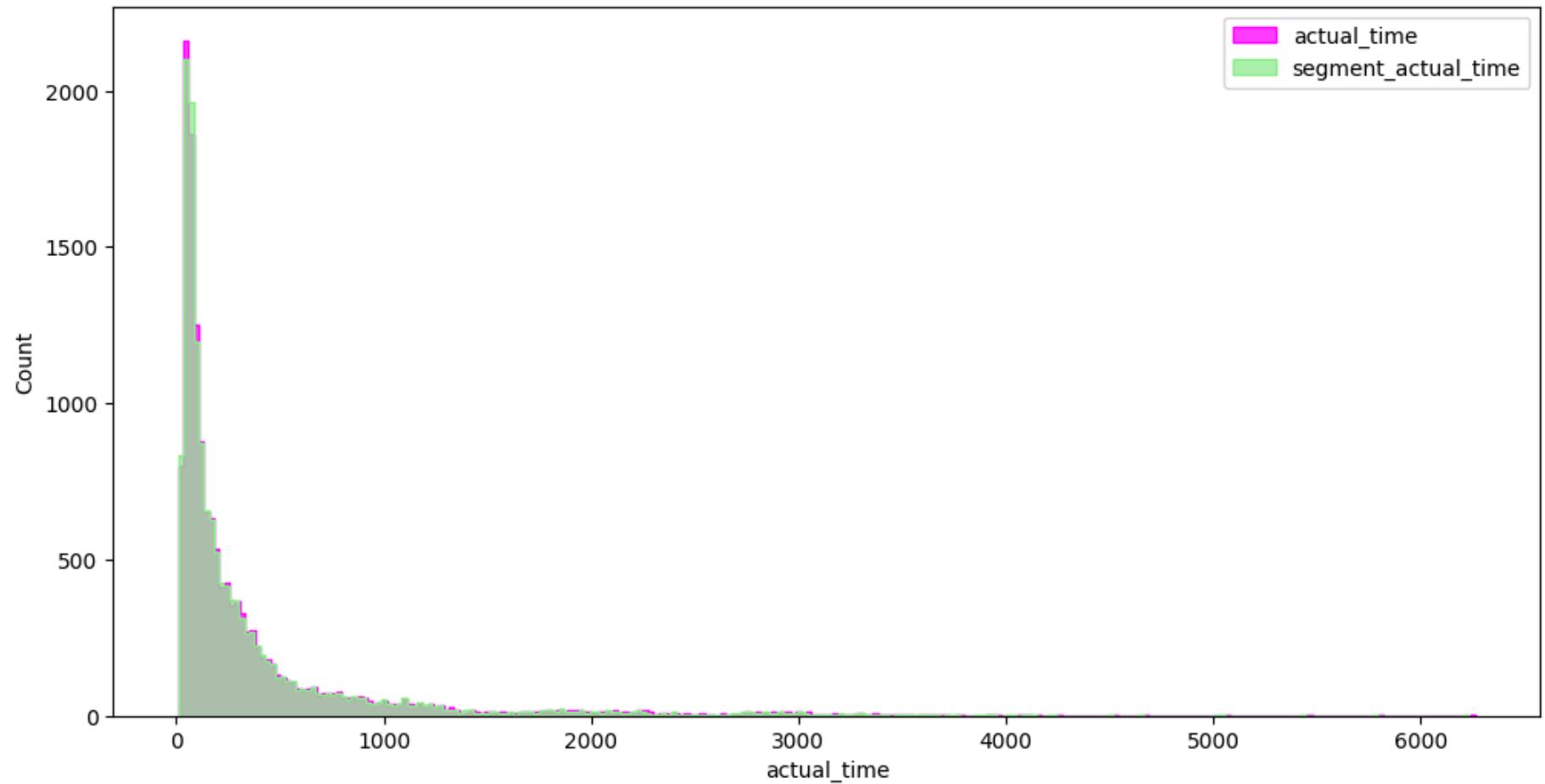
Out[152...

	actual_time	segment_actual_time
count	14817.000000	14817.000000
mean	357.143768	353.892273
std	561.396118	556.247925
min	9.000000	9.000000
25%	67.000000	66.000000
50%	149.000000	147.000000
75%	370.000000	367.000000
max	6265.000000	6230.000000

In [155... `#Visual Tests to know if the samples follow normal distribution`

```
plt.figure(figsize = (12, 6))
sns.histplot(data_2['actual_time'], element = 'step', color = 'magenta')
sns.histplot(data_2['segment_actual_time'], element = 'step', color = 'lightgreen')
plt.legend(['actual_time', 'segment_actual_time'])
plt.plot()
```

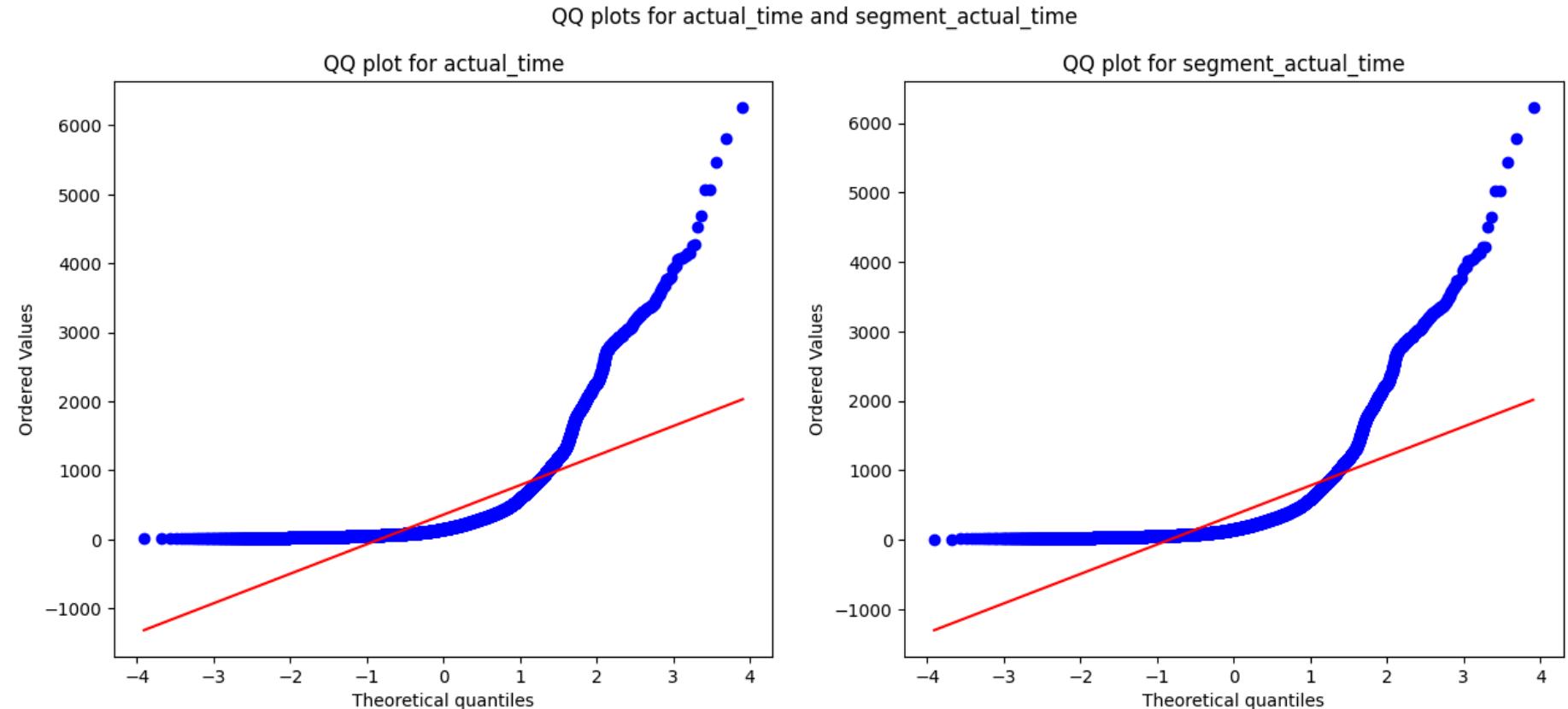
Out[155... `[]`



In [159...]

```
#Distribution check using QQ Plot
plt.figure(figsize = (15, 6))
plt.subplot(1, 2, 1)
plt.suptitle('QQ plots for actual_time and segment_actual_time')
sns.probplot(data_2['actual_time'], plot = plt, dist = 'norm')
plt.title('QQ plot for actual_time')
plt.subplot(1, 2, 2)
sns.probplot(data_2['segment_actual_time'], plot = plt, dist = 'norm')
plt.title('QQ plot for segment_actual_time')
plt.plot()
```

Out[159...]



It can be seen from the above plots that the samples do not come from normal distribution.

Applying Shapiro-Wilk test for normality

H0 : The sample follows normal distribution

H1 : The sample does not follow normal distribution

alpha = 0.05

Test Statistics : Shapiro-Wilk test for normality

```
In [160... test_stat, p_value = spy.shapiro(data_2['actual_time'].sample(5000))
print('p-value', p_value)
if p_value < 0.05:
    print('The sample does not follow normal distribution')
else:
    print('The sample follows normal distribution')
```

```
p-value 0.0
The sample does not follow normal distribution
```

```
In [161... test_stat, p_value = spy.shapiro(data_2['segment_actual_time'].sample(5000))
print('p-value', p_value)
if p_value < 0.05:
    print('The sample does not follow normal distribution')
else:
    print('The sample follows normal distribution')
```

```
p-value 0.0
The sample does not follow normal distribution
```

```
In [162... #Transforming the data using boxcox transformation to check if the transformed data follows normal distribution.
```

```
In [163... transformed_actual_time = spy.boxcox(data_2['actual_time'])[0]
test_stat, p_value = spy.shapiro(transformed_actual_time)
print('p-value', p_value)
if p_value < 0.05:
    print('The sample does not follow normal distribution')
else:
    print('The sample follows normal distribution')
```

```
p-value 1.020620453603145e-28
The sample does not follow normal distribution
```

```
In [164... transformed_segment_actual_time = spy.boxcox(data_2['segment_actual_time'])[0]
test_stat, p_value = spy.shapiro(transformed_segment_actual_time)
print('p-value', p_value)
if p_value < 0.05:
    print('The sample does not follow normal distribution')
else:
    print('The sample follows normal distribution')
```

```
p-value 5.700074948787037e-29  
The sample does not follow normal distribution
```

```
In [165...]: #Even after applying the boxcox transformation on each of the "actual_time" and "segment_actual_time" columns,  
#the distributions do not follow normal distribution.  
  
#Homogeneity of Variances using Lavene's test
```

```
In [167...]: # Null Hypothesis(H0) - Homogenous Variance  
  
# Alternate Hypothesis(HA) - Non Homogenous Variance  
  
test_stat, p_value = spy.levene(data_2['actual_time'], data_2['segment_actual_time'])  
print('p-value', p_value)  
  
if p_value < 0.05:  
    print('The samples do not have Homogenous Variance')  
else:  
    print('The samples have Homogenous Variance ')  
  
p-value 0.695502241317651  
The samples have Homogenous Variance
```

```
In [170...]: #Since the samples do not come from normal distribution T-Test cannot be applied here,  
#we can perform its non parametric equivalent test i.e., Mann-Whitney U rank test for two independent samples.  
test_stat, p_value = spy.mannwhitneyu(data_2['actual_time'], data_2['segment_actual_time'])  
print('p-value', p_value)  
if p_value < 0.05:  
    print('The samples are not similar')  
else:  
    print('The samples are similar ')  
  
p-value 0.4164235159622476  
The samples are similar
```

```
In [171...]: #Since p-value > alpha therfore it can be concluded that actual_time and segment_actual_time are similar.
```

Do hypothesis testing/ visual analysis between osrm distance aggregated value and segment osrm distance aggregated value

(aggregated values are the values you'll get after merging the rows on the basis of trip_uuid)

```
In [173...]: data_2[['osrm_distance', 'segment_osrm_distance']].describe()
```

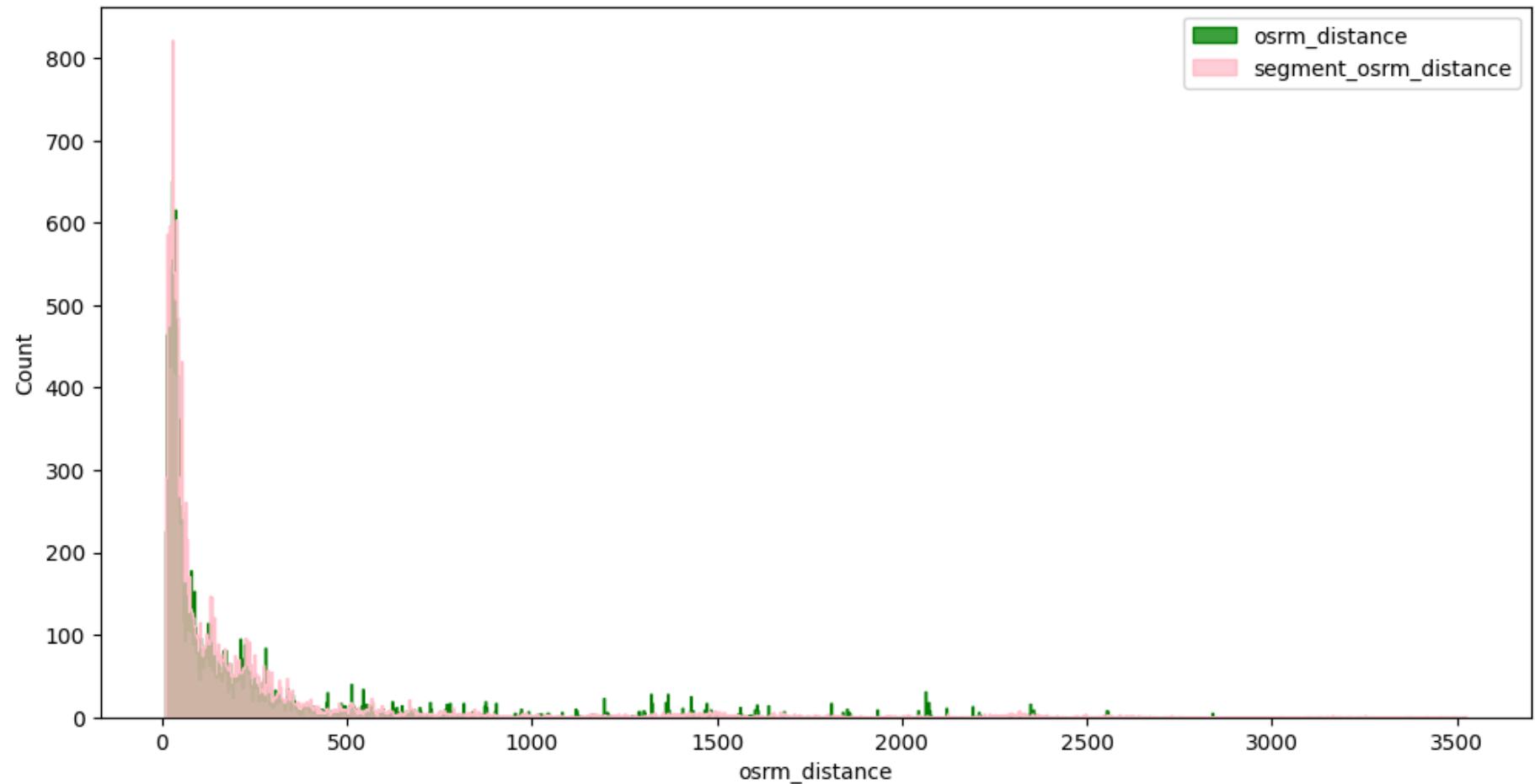
```
Out[173...]:
```

	osrm_distance	segment_osrm_distance
count	14817.000000	14817.000000
mean	204.344711	223.201157
std	370.395569	416.628387
min	9.072900	9.072900
25%	30.819201	32.654499
50%	65.618805	70.154404
75%	208.475006	218.802399
max	2840.081055	3523.632324

```
In [174...]: #Visual Tests to know if the samples follow normal distribution
```

```
In [175...]: plt.figure(figsize = (12, 6))
sns.histplot(data_2['osrm_distance'], element = 'step', color = 'green', bins = 1000)
sns.histplot(data_2['segment_osrm_distance'], element = 'step', color = 'pink', bins = 1000)
plt.legend(['osrm_distance', 'segment_osrm_distance'])
plt.plot()
```

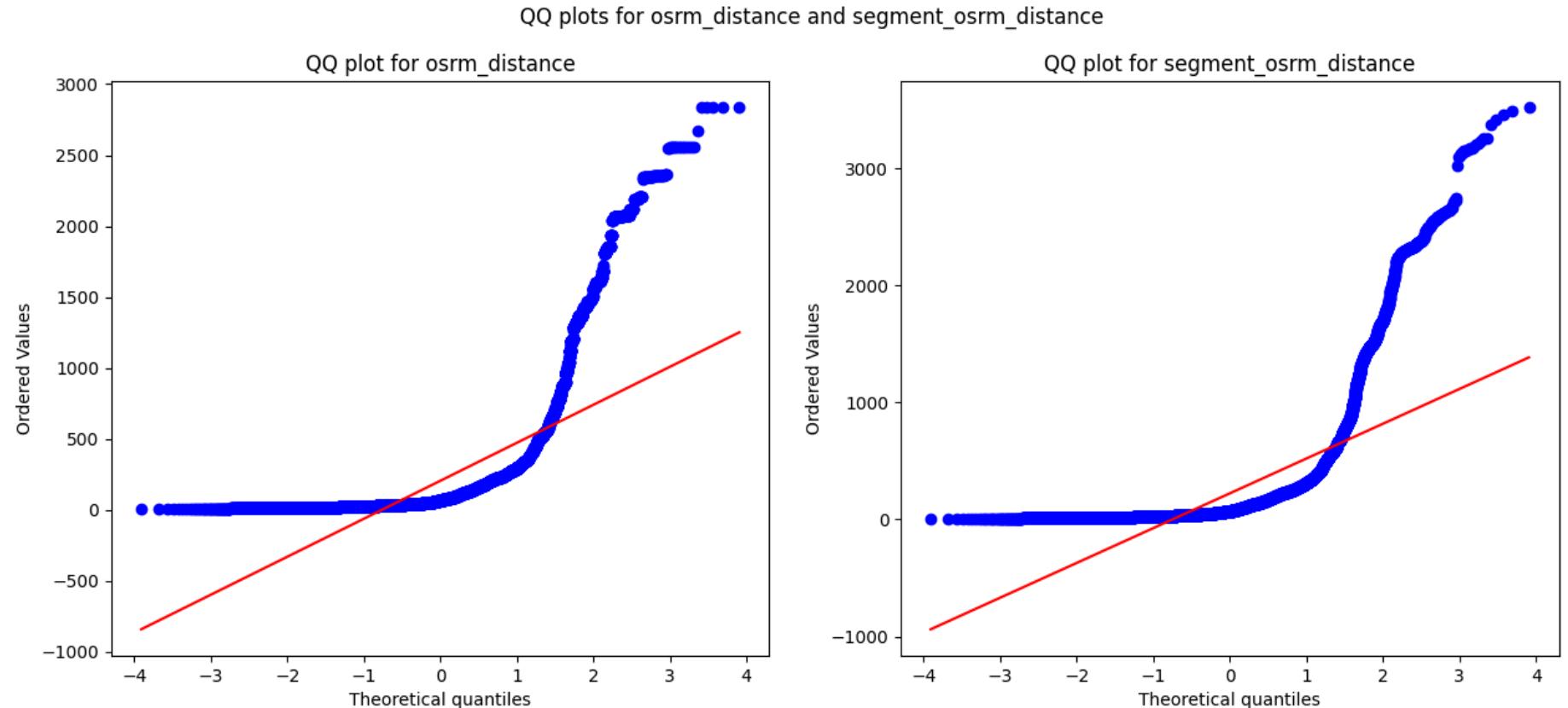
```
Out[175...]: []
```



```
In [176]: #Distribution check using QQ Plot
```

```
In [177]: plt.figure(figsize = (15, 6))
plt.subplot(1, 2, 1)
plt.suptitle('QQ plots for osrm_distance and segment_osrm_distance')
sns.probplot(data_2['osrm_distance'], plot = plt, dist = 'norm')
plt.title('QQ plot for osrm_distance')
plt.subplot(1, 2, 2)
sns.probplot(data_2['segment_osrm_distance'], plot = plt, dist = 'norm')
plt.title('QQ plot for segment_osrm_distance')
plt.plot()
```

Out[177...]



It can be seen from the above plots that the samples do not come from normal distribution.

Applying Shapiro-Wilk test for normality

H0 : The sample follows normal distribution

H1 : The sample does not follow normal distribution

alpha = 0.05

Test Statistics : Shapiro-Wilk test for normality

```
In [178... test_stat, p_value = spy.shapiro(data_2['osrm_distance'].sample(5000))
print('p-value', p_value)
if p_value < 0.05:
    print('The sample does not follow normal distribution')
else:
    print('The sample follows normal distribution')
```

```
p-value 0.0
The sample does not follow normal distribution
```

```
In [179... test_stat, p_value = spy.shapiro(data_2['segment_osrm_distance'].sample(5000))
print('p-value', p_value)
if p_value < 0.05:
    print('The sample does not follow normal distribution')
else:
    print('The sample follows normal distribution')
```

```
p-value 0.0
The sample does not follow normal distribution
```

```
In [180... #Transforming the data using boxcox transformation to check if the transformed data follows normal distribution.
```

```
In [181... transformed_osrm_distance = spy.boxcox(data_2['osrm_distance'])[0]
test_stat, p_value = spy.shapiro(transformed_osrm_distance)
print('p-value', p_value)
if p_value < 0.05:
    print('The sample does not follow normal distribution')
else:
    print('The sample follows normal distribution')
```

```
p-value 7.063104779582808e-41
The sample does not follow normal distribution
```

```
In [182... transformed_segment_osrm_distance = spy.boxcox(data_2['segment_osrm_distance'])[0]
test_stat, p_value = spy.shapiro(transformed_segment_osrm_distance)
print('p-value', p_value)
if p_value < 0.05:
    print('The sample does not follow normal distribution')
else:
    print('The sample follows normal distribution')
```

```
p-value 3.049169406432229e-38  
The sample does not follow normal distribution
```

```
In [183...]: #Even after applying the boxcox transformation on each of the "osrm_distance" and "segment_osrm_distance" columns,  
#the distributions do not follow normal distribution.
```

```
#Homogeneity of Variances using Lavene's test
```

```
In [184...]: # Null Hypothesis(H0) - Homogenous Variance  
  
# Alternate Hypothesis(HA) - Non Homogenous Variance  
  
test_stat, p_value = spy.levene(data_2['osrm_distance'], data_2['segment_osrm_distance'])  
print('p-value', p_value)  
  
if p_value < 0.05:  
    print('The samples do not have Homogenous Variance')  
else:  
    print('The samples have Homogenous Variance ')
```

```
p-value 0.00020976006524780905  
The samples do not have Homogenous Variance
```

```
In [185...]: #Since the samples do not follow any of the assumptions, T-Test cannot be applied here.  
#We can perform its non parametric equivalent test i.e., Mann-Whitney U rank test for two independent samples.
```

```
In [186...]: test_stat, p_value = spy.mannwhitneyu(data_2['osrm_distance'], data_2['segment_osrm_distance'])  
print('p-value', p_value)  
if p_value < 0.05:  
    print('The samples are not similar')  
else:  
    print('The samples are similar ')
```

```
p-value 9.509410818847664e-07  
The samples are not similar
```

```
In [187...]: #Since p-value < alpha therfore it can be concluded that osrm_distance and segment_osrm_distance are not similar.
```

**Do hypothesis testing/ visual analysis between osrm time aggregated value and segment osrm time aggregated value
(aggregated values are the values you'll get after merging the rows on the basis of trip_uuid)**

In [188... `data_2[['osrm_time', 'segment_osrm_time']].describe().T`

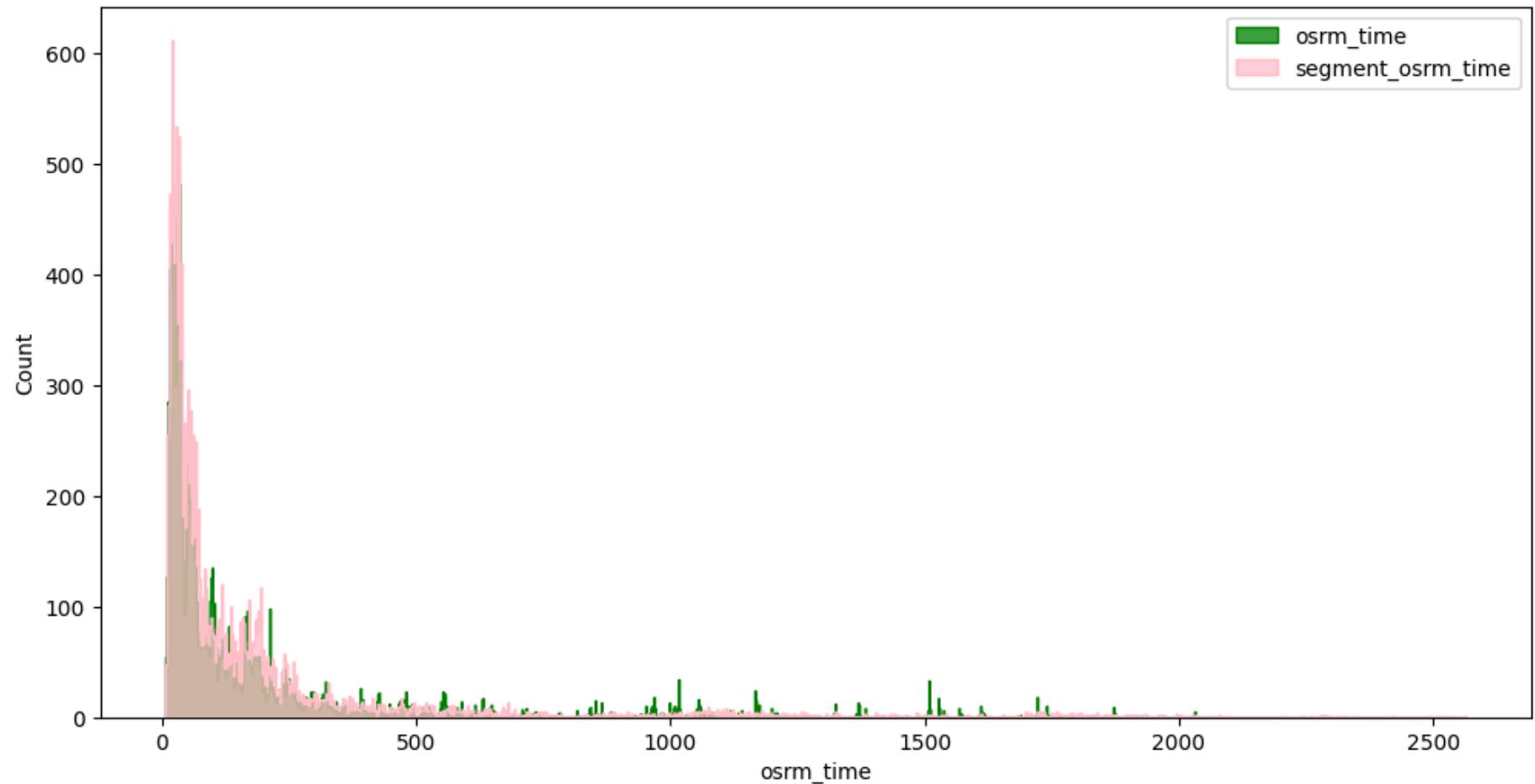
Out[188...

	count	mean	std	min	25%	50%	75%	max
osrm_time	14817.0	161.384018	271.360992	6.0	29.0	60.0	168.0	2032.0
segment_osrm_time	14817.0	180.949783	314.542053	6.0	31.0	65.0	185.0	2564.0

In [189... `#Visual Tests to know if the samples follow normal distribution`

In [190... `plt.figure(figsize = (12, 6))
sns.histplot(data_2['osrm_time'], element = 'step', color = 'green', bins = 1000)
sns.histplot(data_2['segment_osrm_time'], element = 'step', color = 'pink', bins = 1000)
plt.legend(['osrm_time', 'segment_osrm_time'])
plt.plot()`

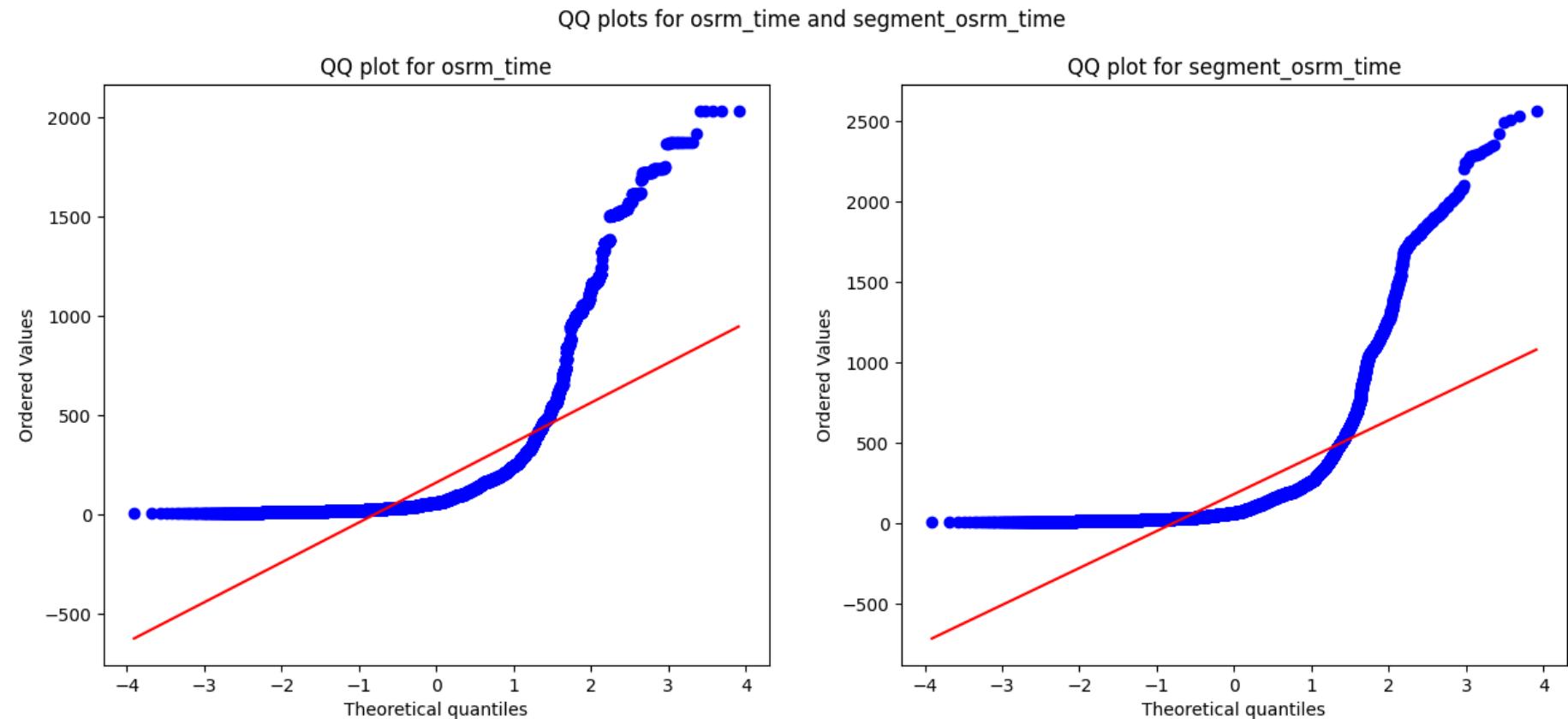
Out[190... `[]`



```
In [191]: #Distribution check using QQ Plot
```

```
In [192]: plt.figure(figsize = (15, 6))
plt.subplot(1, 2, 1)
plt.suptitle('QQ plots for osrm_time and segment_osrm_time')
sns.probplot(data_2['osrm_time'], plot = plt, dist = 'norm')
plt.title('QQ plot for osrm_time')
plt.subplot(1, 2, 2)
sns.probplot(data_2['segment_osrm_time'], plot = plt, dist = 'norm')
plt.title('QQ plot for segment_osrm_time')
plt.plot()
```

Out[192...]



It can be seen from the above plots that the samples do not come from normal distribution.

Applying Shapiro-Wilk test for normality

H0 : The sample follows normal distribution

H1 : The sample does not follow normal distribution

alpha = 0.05

Test Statistics : Shapiro-Wilk test for normality

```
In [193...  
    test_stat, p_value = spy.shapiro(data_2['osrm_time'].sample(5000))  
    print('p-value', p_value)  
    if p_value < 0.05:  
        print('The sample does not follow normal distribution')  
    else:  
        print('The sample follows normal distribution')
```

```
p-value 0.0  
The sample does not follow normal distribution
```

```
In [194...  
    test_stat, p_value = spy.shapiro(data_2['segment_osrm_time'].sample(5000))  
    print('p-value', p_value)  
    if p_value < 0.05:  
        print('The sample does not follow normal distribution')  
    else:  
        print('The sample follows normal distribution')
```

```
p-value 0.0  
The sample does not follow normal distribution
```

```
In [195... #Transforming the data using boxcox transformation to check if the transformed data follows normal distribution.
```

```
In [196...  
    transformed_osrm_time = spy.boxcox(data_2['osrm_time'])[0]  
    test_stat, p_value = spy.shapiro(transformed_osrm_time)  
    print('p-value', p_value)  
    if p_value < 0.05:  
        print('The sample does not follow normal distribution')  
    else:  
        print('The sample follows normal distribution')
```

```
p-value 3.5882550510138333e-35  
The sample does not follow normal distribution
```

```
In [197...  
    transformed_segment_osrm_time = spy.boxcox(data_2['segment_osrm_time'])[0]  
    test_stat, p_value = spy.shapiro(transformed_segment_osrm_time)  
    print('p-value', p_value)  
    if p_value < 0.05:  
        print('The sample does not follow normal distribution')  
    else:  
        print('The sample follows normal distribution')
```

```
p-value 4.943039152219146e-34  
The sample does not follow normal distribution
```

```
In [198...]: #Even after applying the boxcox transformation on each of the "osrm_time" and "segment_osrm_time" columns,  
#the distributions do not follow normal distribution.  
  
#Homogeneity of Variances using Lavene's test
```

```
In [199...]: # Null Hypothesis(H0) - Homogenous Variance  
  
# Alternate Hypothesis(HA) - Non Homogenous Variance  
  
test_stat, p_value = spy.levene(data_2['osrm_time'], data_2['segment_osrm_time'])  
print('p-value', p_value)  
  
if p_value < 0.05:  
    print('The samples do not have Homogenous Variance')  
else:  
    print('The samples have Homogenous Variance ')  
  
p-value 8.349506135727595e-08  
The samples do not have Homogenous Variance
```

```
In [201...]: #Since the samples do not follow any of the assumptions, T-Test cannot be applied here.  
#We can perform its non parametric equivalent test i.e., Mann-Whitney U rank test for two independent samples.
```

```
In [202...]: test_stat, p_value = spy.mannwhitneyu(data_2['osrm_time'], data_2['segment_osrm_time'])  
print('p-value', p_value)  
if p_value < 0.05:  
    print('The samples are not similar')  
else:  
    print('The samples are similar ')  
  
p-value 2.2995370859748865e-08  
The samples are not similar
```

```
In [203...]: #Since p-value < alpha therfore it can be concluded that osrm_time and segment_osrm_time are not similar.
```

Find outliers in the numerical variables (you might find outliers in almost all the variables), and check it using visual analysis

In [204...]

```
numerical_columns = ['od_total_time', 'start_scan_to_end_scan', 'actual_distance_to_destination',
                     'actual_time', 'osrm_time', 'osrm_distance', 'segment_actual_time',
                     'segment_osrm_time', 'segment_osrm_distance']
data_2[numerical_columns].describe().T
```

Out[204...]

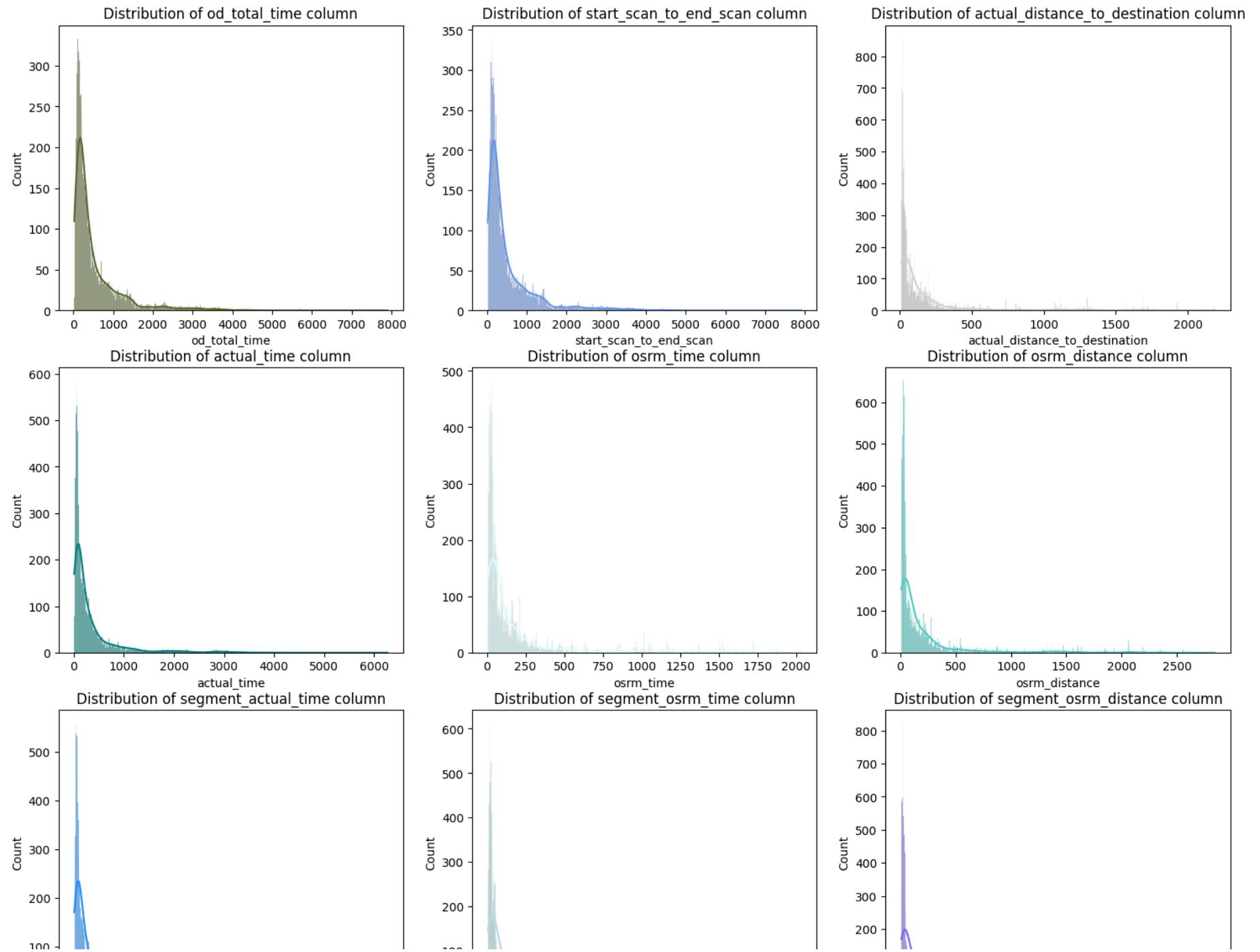
		count	mean	std	min	25%	50%	75%	max
	od_total_time	14817.0	531.697630	658.868223	23.460000	149.930000	280.770000	638.200000	7898.550000
	start_scan_to_end_scan	14817.0	530.810016	658.705957	23.000000	149.000000	280.000000	637.000000	7898.000000
	actual_distance_to_destination	14817.0	164.477829	305.388153	9.002461	22.837238	48.474072	164.583206	2186.531738
	actual_time	14817.0	357.143768	561.396118	9.000000	67.000000	149.000000	370.000000	6265.000000
	osrm_time	14817.0	161.384018	271.360992	6.000000	29.000000	60.000000	168.000000	2032.000000
	osrm_distance	14817.0	204.344711	370.395569	9.072900	30.819201	65.618805	208.475006	2840.081055
	segment_actual_time	14817.0	353.892273	556.247925	9.000000	66.000000	147.000000	367.000000	6230.000000
	segment_osrm_time	14817.0	180.949783	314.542053	6.000000	31.000000	65.000000	185.000000	2564.000000
	segment_osrm_distance	14817.0	223.201157	416.628387	9.072900	32.654499	70.154404	218.802399	3523.632324

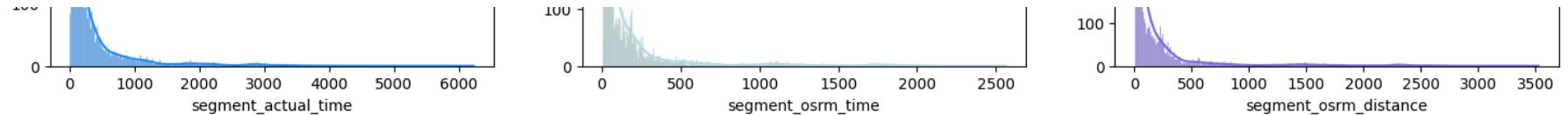
In [206...]

```
import matplotlib as mpl
```

In [207...]

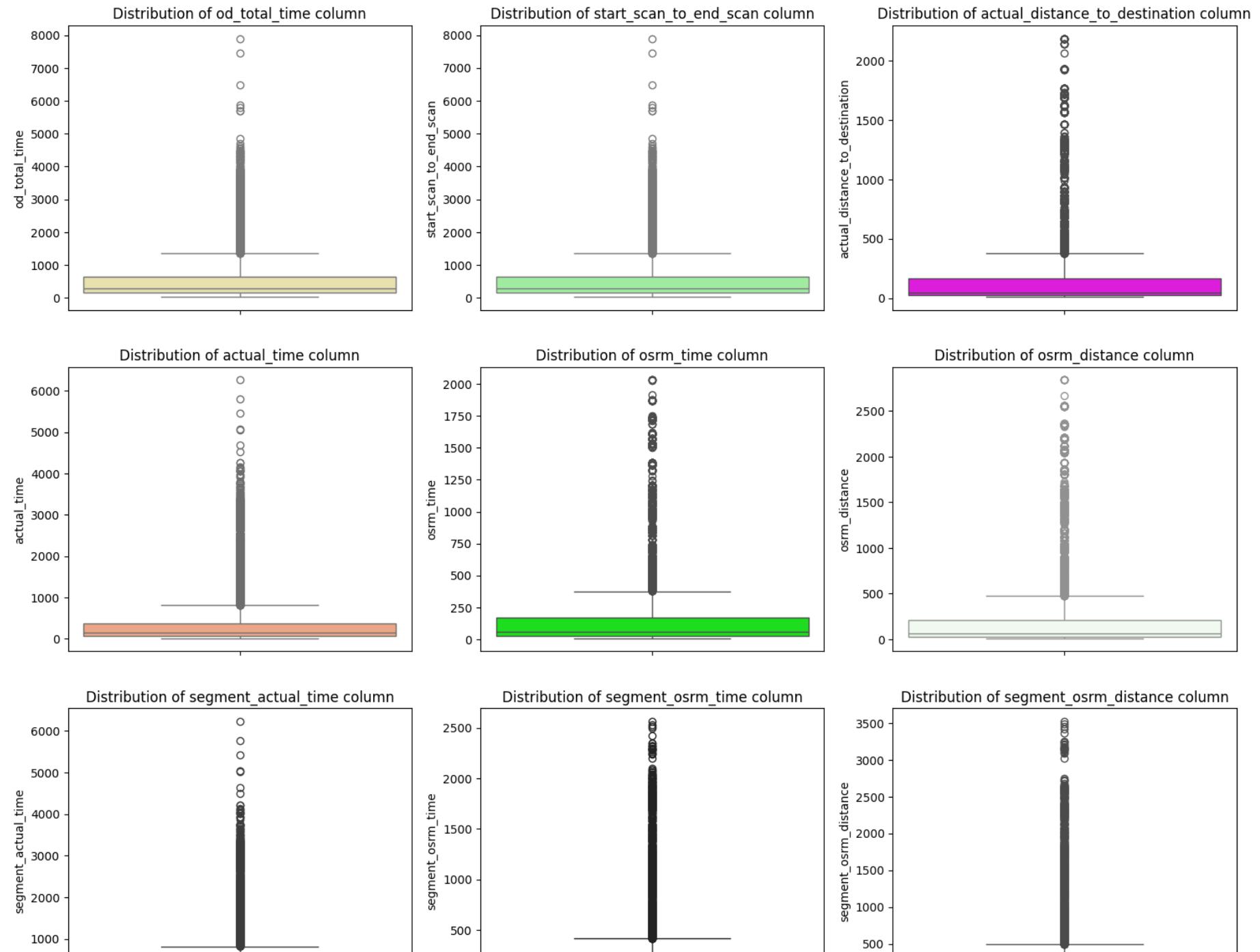
```
plt.figure(figsize = (18, 15))
for i in range(len(numerical_columns)):
    plt.subplot(3, 3, i + 1)
    clr = np.random.choice(list(mpl.colors.cnames))
    sns.histplot(data_2[numerical_columns[i]], bins = 1000, kde = True, color = clr)
    plt.title(f"Distribution of {numerical_columns[i]} column")
    plt.plot()
```

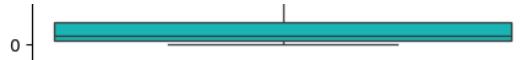




```
In [209]: #It can be inferred from the above plots that data in all the numerical columns are right skewed
```

```
In [210]: plt.figure(figsize = (18, 15))
for i in range(len(numerical_columns)):
    plt.subplot(3, 3, i + 1)
    clr = np.random.choice(list(mpl.colors.cnames))
    sns.boxplot(data_2[numerical_columns[i]], color = clr)
    plt.title(f"Distribution of {numerical_columns[i]} column")
    plt.plot()
```





```
In [211... #It can be clearly seen in the above plots that there are outliers in all the numerical columns that need to be treated.
```

```
In [214... # Detecting Outliers
```

```
for i in numerical_columns:  
    Q1 = np.quantile(data_2[i], 0.25)  
    Q3 = np.quantile(data_2[i], 0.75)  
    IQR = Q3 - Q1  
    LB = Q1 - 1.5 * IQR  
    UB = Q3 + 1.5 * IQR  
    outliers = data_2.loc[(data_2[i] < LB) | (data_2[i] > UB)]  
    print('Column :', i)  
    print(f'Q1 : {Q1}')  
    print(f'Q3 : {Q3}')  
    print(f'IQR : {IQR}')  
    print(f'LB : {LB}')  
    print(f'UB : {UB}')  
    print(f'Number of outliers : {outliers.shape[0]}')  
    print('-----')
```

Column : od_total_time

Q1 : 149.93

Q3 : 638.2

IQR : 488.2700000000004

LB : -582.475000000001

UB : 1370.605

Number of outliers : 1266

Column : start_scan_to_end_scan

Q1 : 149.0

Q3 : 637.0

IQR : 488.0

LB : -583.0

UB : 1369.0

Number of outliers : 1267

Column : actual_distance_to_destination

Q1 : 22.837238311767578

Q3 : 164.5832061767578

IQR : 141.74596786499023

LB : -189.78171348571777

UB : 377.20215797424316

Number of outliers : 1449

Column : actual_time

Q1 : 67.0

Q3 : 370.0

IQR : 303.0

LB : -387.5

UB : 824.5

Number of outliers : 1643

Column : osrm_time

Q1 : 29.0

Q3 : 168.0

IQR : 139.0

LB : -179.5

UB : 376.5

Number of outliers : 1517

Column : osrm_distance

```
Q1 : 30.81920051574707
Q3 : 208.47500610351562
IQR : 177.65580558776855
LB : -235.66450786590576
UB : 474.95871448516846
Number of outliers : 1524
-----
Column : segment_actual_time
Q1 : 66.0
Q3 : 367.0
IQR : 301.0
LB : -385.5
UB : 818.5
Number of outliers : 1643
-----
Column : segment_osrm_time
Q1 : 31.0
Q3 : 185.0
IQR : 154.0
LB : -200.0
UB : 416.0
Number of outliers : 1492
-----
Column : segment_osrm_distance
Q1 : 32.65449905395508
Q3 : 218.80239868164062
IQR : 186.14789962768555
LB : -246.56735038757324
UB : 498.02424812316895
Number of outliers : 1548
```

In [215...]

```
#The outliers present in our sample data can be the true outliers. It's best to remove outliers only when there is a sound reason.
#Some outliers represent natural variations in the population, and they should be left as is in the dataset.
```

Do one-hot encoding of categorical variables (like route_type)

```
In [216...]: # Get value counts before one-hot encoding  
  
data_2['route_type'].value_counts()
```

```
Out[216...]: route_type  
Carting      8908  
FTL         5909  
Name: count, dtype: int64
```

```
In [217...]: # Perform one-hot encoding on categorical column route type
```

```
In [218...]: from sklearn.preprocessing import LabelEncoder  
label_encoder = LabelEncoder()  
data_2['route_type'] = label_encoder.fit_transform(data_2['route_type'])
```

```
In [219...]: # Get value counts after one-hot encoding  
  
data_2['route_type'].value_counts()
```

```
Out[219...]: route_type  
0      8908  
1      5909  
Name: count, dtype: int64
```

```
In [220...]: # Get value counts of categorical variable 'data' before one-hot encoding  
  
data_2['data'].value_counts()
```

```
Out[220...]: data  
training    10654  
test        4163  
Name: count, dtype: int64
```

```
In [221...]: # Perform one-hot encoding on categorical variable 'data'
```

```
In [222...]: label_encoder = LabelEncoder()  
data_2['data'] = label_encoder.fit_transform(data_2['data'])
```

```
In [223...]: # Get value counts after one-hot encoding  
  
data_2['data'].value_counts()
```

```
Out[223...]: data  
1    10654  
0     4163  
Name: count, dtype: int64
```

Normalize/ Standardize the numerical features using MinMaxScaler or StandardScaler.

```
In [224...]: from sklearn.preprocessing import MinMaxScaler
```

```
In [225...]: plt.figure(figsize = (10, 6))  
scaler = MinMaxScaler()  
scaled = scaler.fit_transform(data_2['od_total_time'].to_numpy().reshape(-1, 1))  
sns.histplot(scaled)  
plt.title(f"Normalized {data_2['od_total_time']} column")  
plt.legend('od_total_time')  
plt.plot()
```

```
Out[225...]: []
```

Normalized 0 2260.11

1 181.61
2 3934.36
3 100.49
4 718.34

...
14812 258.03
14813 60.59
14814 422.12
14815 348.52
14816 354.40

Name: od_total_time, Length: 14817, dtype: float64 column





In [226...]

```
plt.figure(figsize = (10, 6))
scaler = MinMaxScaler()
scaled = scaler.fit_transform(data_2['start_scan_to_end_scan'].to_numpy().reshape(-1, 1))
sns.histplot(scaled)
plt.title(f"Normalized {data_2['start_scan_to_end_scan']} column")
plt.plot()
```

Out[226...]

[]

Normalized 0 2259.0

1 180.0

2 3933.0

3 100.0

4 717.0

...

14812 257.0

14813 60.0

14814 421.0

14815 347.0

14816 353.0

Name: start_scan_to_end_scan, Length: 14817, dtype: float64 column





```
In [227]: plt.figure(figsize = (10, 6))
scaler = MinMaxScaler()
scaled = scaler.fit_transform(data_2['actual_distance_to_destination'].to_numpy().reshape(-1, 1))
sns.histplot(scaled)
plt.title(f"Normalized {data_2['actual_distance_to_destination']} column")
plt.plot()
```

```
Out[227]: []
```

```
Normalized 0      824.732849
           1      73.186905
           2     1927.404297
           3     17.175274
           4    127.448502
...
14812    57.762333
14813   15.513784
14814   38.684837
14815  134.723831
14816   66.081528
```

Name: actual_distance_to_destination, Length: 14817, dtype: float32 column





In [228]:

```
plt.figure(figsize = (10, 6))
scaler = MinMaxScaler()
scaled = scaler.fit_transform(data_2['actual_time'].to_numpy().reshape(-1, 1))
sns.histplot(scaled)
plt.title(f"Normalized {data_2['actual_time']} column")
plt.plot()
```

Out[228]: []

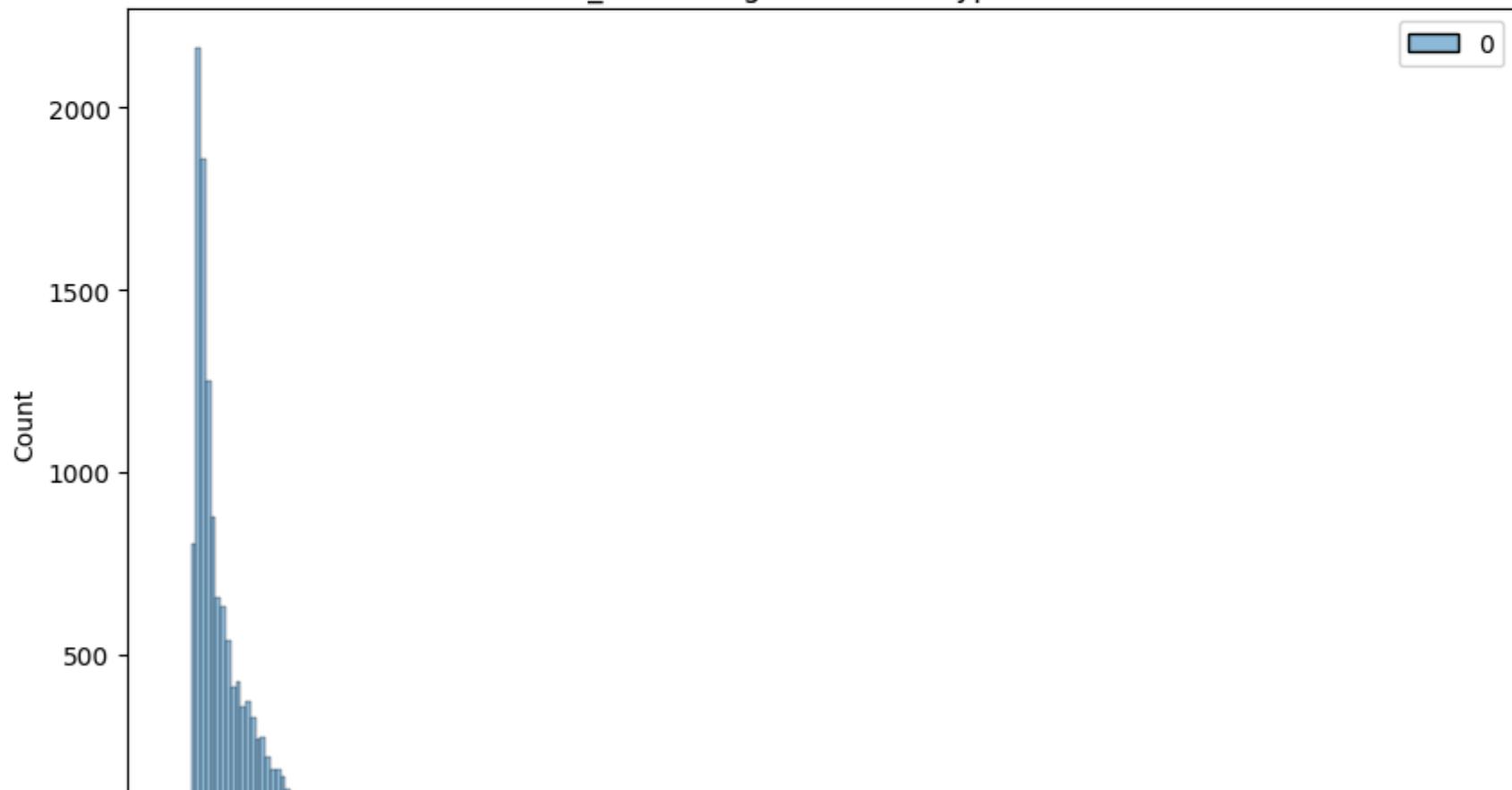
Normalized 0 1562.0

1 143.0
2 3347.0
3 59.0
4 341.0

...

14812 83.0
14813 21.0
14814 282.0
14815 264.0
14816 275.0

Name: actual_time, Length: 14817, dtype: float32 column





In [229...]

```
plt.figure(figsize = (10, 6))
scaler = MinMaxScaler()
scaled = scaler.fit_transform(data_2['osrm_time'].to_numpy().reshape(-1, 1))
sns.histplot(scaled)
plt.title(f"Normalized {data_2['osrm_time']} column")
plt.plot()
```

Out[229...]

[]

Normalized 0 717.0

1 68.0
2 1740.0
3 15.0
4 117.0

...

14812 62.0
14813 12.0
14814 48.0
14815 179.0
14816 68.0

Name: osrm_time, Length: 14817, dtype: float32 column





In [230...]

```
plt.figure(figsize = (10, 6))
scaler = MinMaxScaler()
scaled = scaler.fit_transform(data_2['osrm_distance'].to_numpy().reshape(-1, 1))
sns.histplot(scaled)
plt.title(f"Normalized {data_2['osrm_distance']} column")
plt.plot()
```

Out[230...]

[]

Normalized 0 991.352295

1 85.111000

2 2354.066650

3 19.680000

4 146.791794

...

14812 73.462997

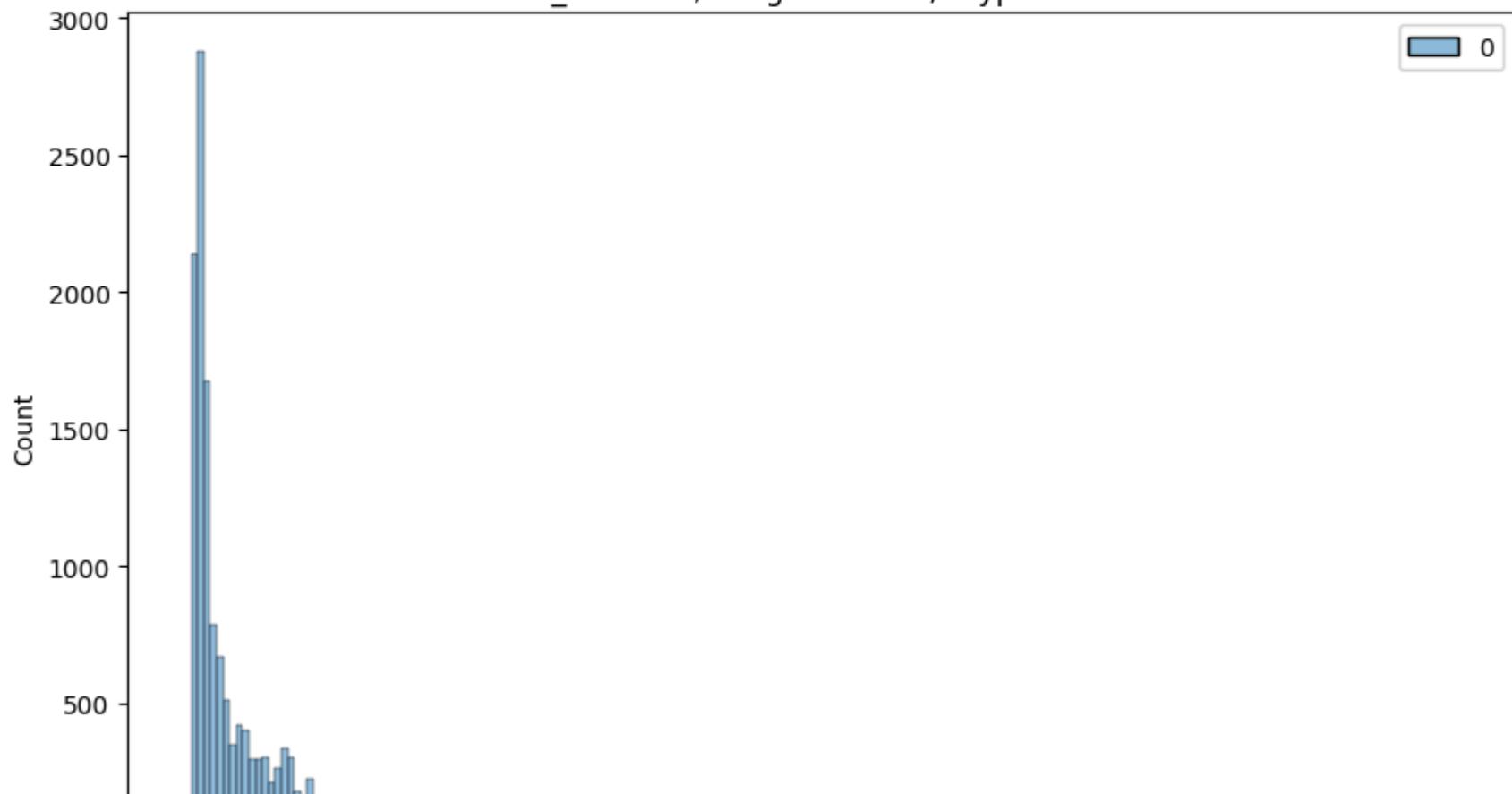
14813 16.088200

14814 58.903702

14815 171.110306

14816 80.578705

Name: osrm_distance, Length: 14817, dtype: float32 column





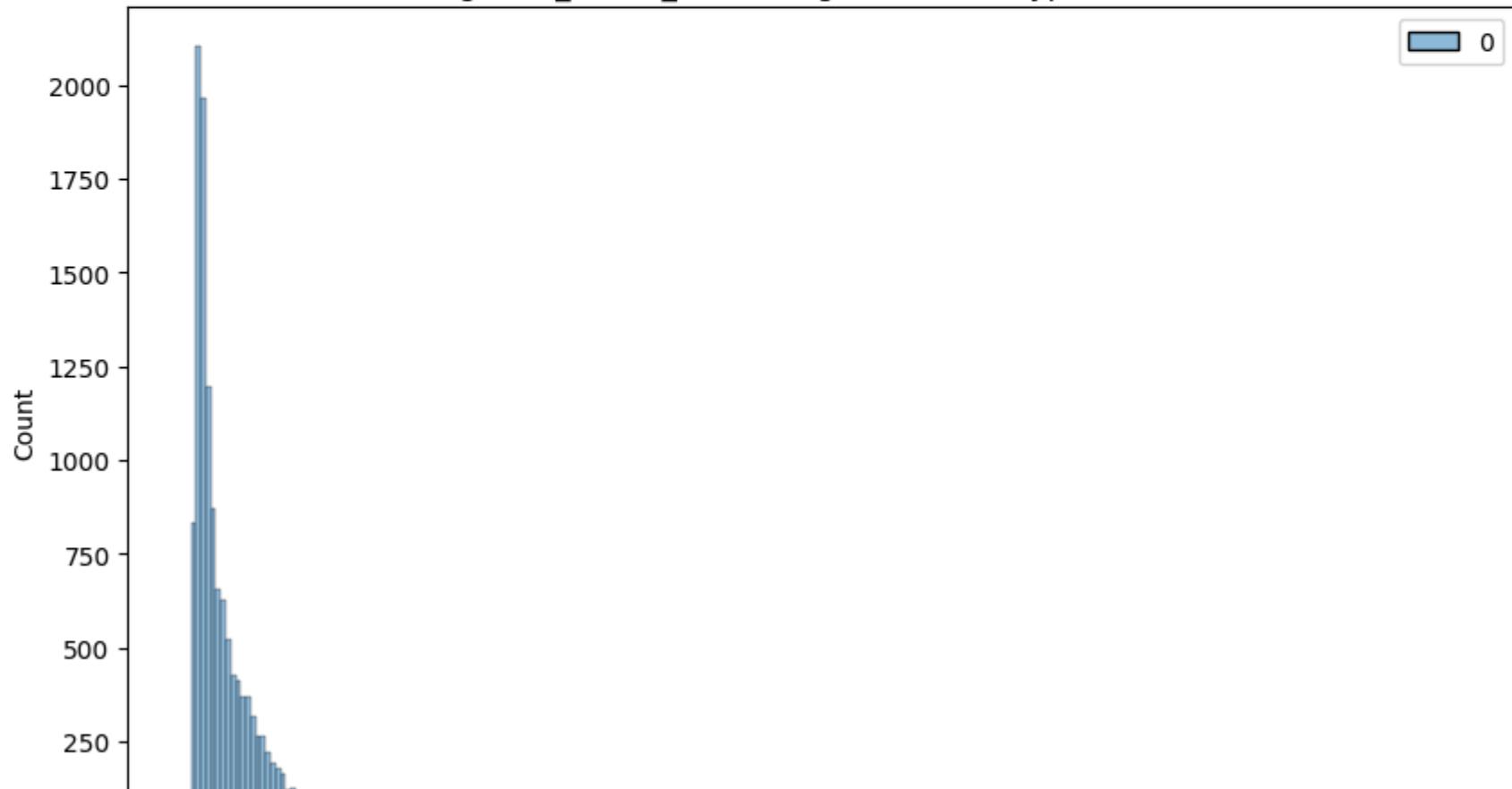
In [231...]

```
plt.figure(figsize = (10, 6))
scaler = MinMaxScaler()
scaled = scaler.fit_transform(data_2['segment_actual_time'].to_numpy().reshape(-1, 1))
sns.histplot(scaled)
plt.title(f"Normalized {data_2['segment_actual_time']} column")
plt.plot()
```

Out[231...]

[]

```
Normalized 0      1548.0
           1      141.0
           2     3308.0
           3      59.0
           4     340.0
...
           14812    82.0
           14813   21.0
           14814   281.0
           14815   258.0
           14816   274.0
Name: segment_actual_time, Length: 14817, dtype: float32 column
```





In [232...]

```
plt.figure(figsize = (10, 6))
scaler = MinMaxScaler()
scaled = scaler.fit_transform(data_2['segment_osrm_time'].to_numpy().reshape(-1, 1))
sns.histplot(scaled)
plt.title(f"Normalized {data_2['segment_osrm_time']} column")
plt.plot()
```

Out[232...]

[]

Normalized 0 1008.0

1 65.0

2 1941.0

3 16.0

4 115.0

...

14812 62.0

14813 11.0

14814 88.0

14815 221.0

14816 67.0

Name: segment_osrm_time, Length: 14817, dtype: float32 column





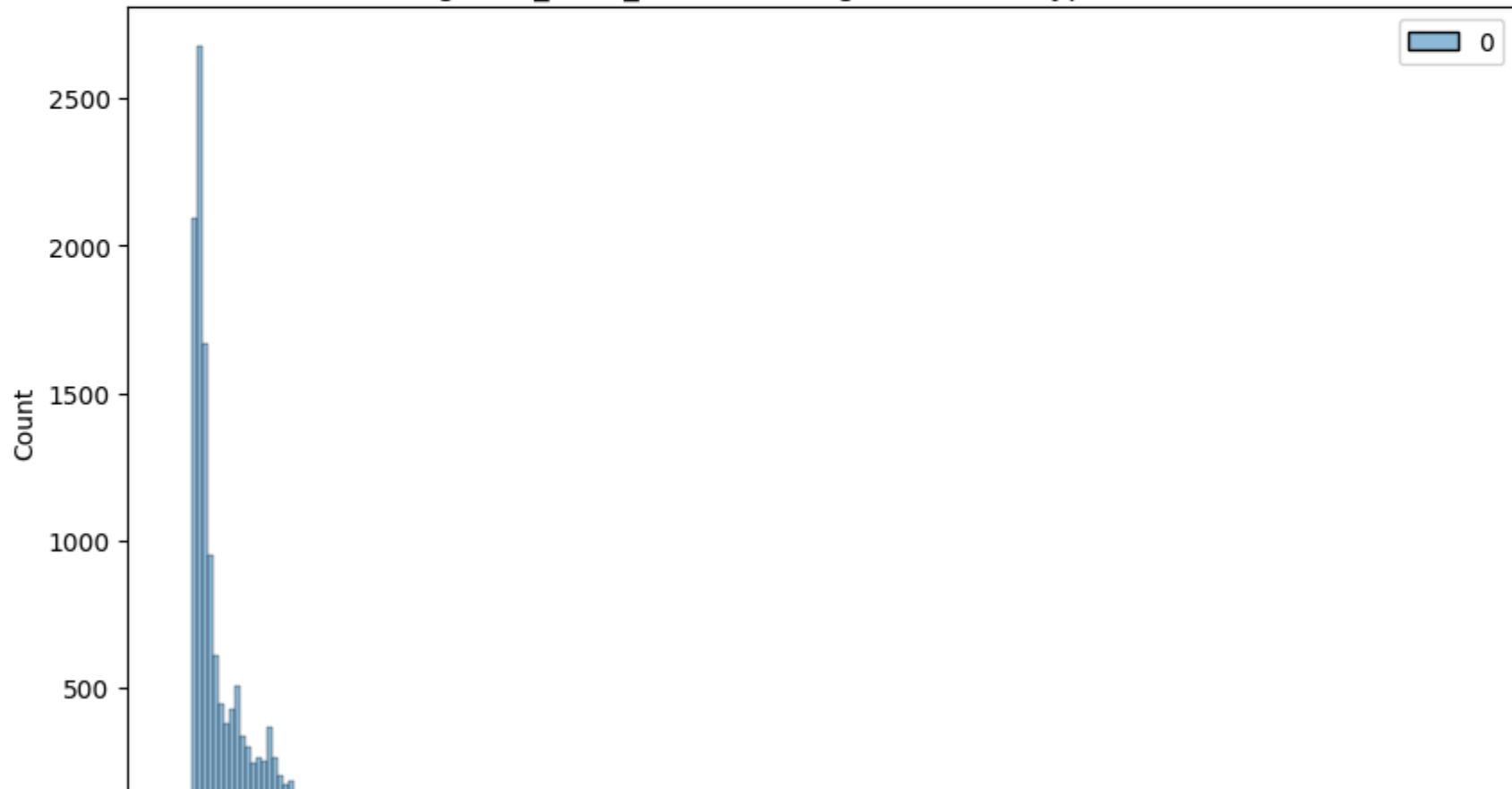
In [233...]

```
plt.figure(figsize = (10, 6))
scaler = MinMaxScaler()
scaled = scaler.fit_transform(data_2['segment_osrm_distance'].to_numpy().reshape(-1, 1))
sns.histplot(scaled)
plt.title(f"Normalized {data_2['segment_osrm_distance']} column")
plt.plot()
```

Out[233...]

[]

```
Normalized 0      1320.473267
              1      84.189400
              2     2545.267822
              3     19.876600
              4    146.791901
              ...
              14812   64.855103
              14813   16.088299
              14814   104.886597
              14815   223.532394
              14816   80.578705
Name: segment_osrm_distance, Length: 14817, dtype: float32 column
```





Column Standardization

```
In [234...]: from sklearn.preprocessing import StandardScaler
```

```
In [235...]: plt.figure(figsize = (10, 6))
# define standard scaler
scaler = StandardScaler()
# transform data
scaled = scaler.fit_transform(data_2['od_total_time'].to_numpy().reshape(-1, 1))
sns.histplot(scaled)
plt.title(f"Standardized {data_2['od_total_time']} column")
plt.legend('od_total_time')
plt.plot()
```

```
Out[235...]: []
```

Standardized 0 2260.11

1 181.61
2 3934.36
3 100.49
4 718.34

...
14812 258.03
14813 60.59
14814 422.12
14815 348.52
14816 354.40

Name: od_total_time, Length: 14817, dtype: float64 column





In [236...]

```
plt.figure(figsize = (10, 6))
scaler = StandardScaler()
scaled = scaler.fit_transform(data_2['start_scan_to_end_scan'].to_numpy().reshape(-1, 1))
sns.histplot(scaled)
plt.title(f"Standardized {data_2['start_scan_to_end_scan']} column")
plt.plot()
```

Out[236...]

[]

Standardized 0 2259.0

1 180.0

2 3933.0

3 100.0

4 717.0

...

14812 257.0

14813 60.0

14814 421.0

14815 347.0

14816 353.0

Name: start_scan_to_end_scan, Length: 14817, dtype: float64 column





```
In [237]: plt.figure(figsize = (10, 6))
scaler = StandardScaler()
scaled = scaler.fit_transform(data_2['actual_distance_to_destination'].to_numpy().reshape(-1, 1))
sns.histplot(scaled)
plt.title(f"Standardized {data_2['actual_distance_to_destination']} column")
plt.plot()
```

```
Out[237]: []
```

Standardized 0 824.732849

1 73.186905

2 1927.404297

3 17.175274

4 127.448502

...

14812 57.762333

14813 15.513784

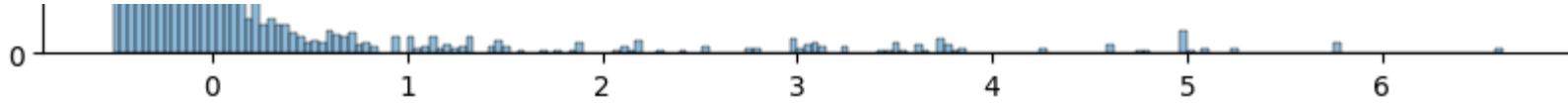
14814 38.684837

14815 134.723831

14816 66.081528

Name: actual_distance_to_destination, Length: 14817, dtype: float32 column





In [238...]

```
plt.figure(figsize = (10, 6))
scaler = StandardScaler()
scaled = scaler.fit_transform(data_2['actual_time'].to_numpy().reshape(-1, 1))
sns.histplot(scaled)
plt.title(f"Standardized {data_2['actual_time']} column")
plt.plot()
```

Out[238...]

[]

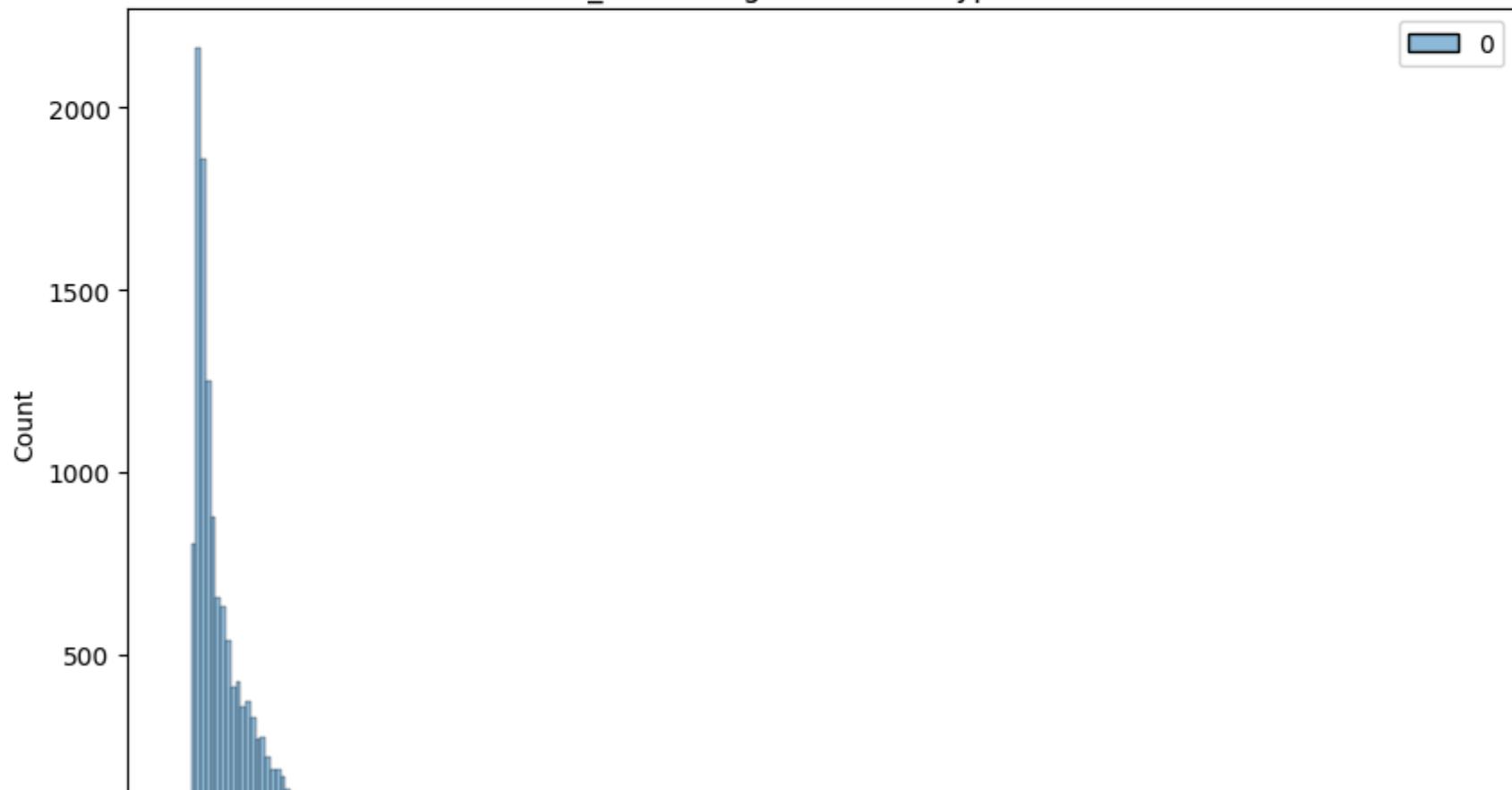
Standardized 0 1562.0

1 143.0
2 3347.0
3 59.0
4 341.0

...

14812 83.0
14813 21.0
14814 282.0
14815 264.0
14816 275.0

Name: actual_time, Length: 14817, dtype: float32 column





In [239...]

```
plt.figure(figsize = (10, 6))
scaler = StandardScaler()
scaled = scaler.fit_transform(data_2['osrm_time'].to_numpy().reshape(-1, 1))
sns.histplot(scaled)
plt.title(f"Standardized {data_2['osrm_time']} column")
plt.plot()
```

Out[239...]

[]

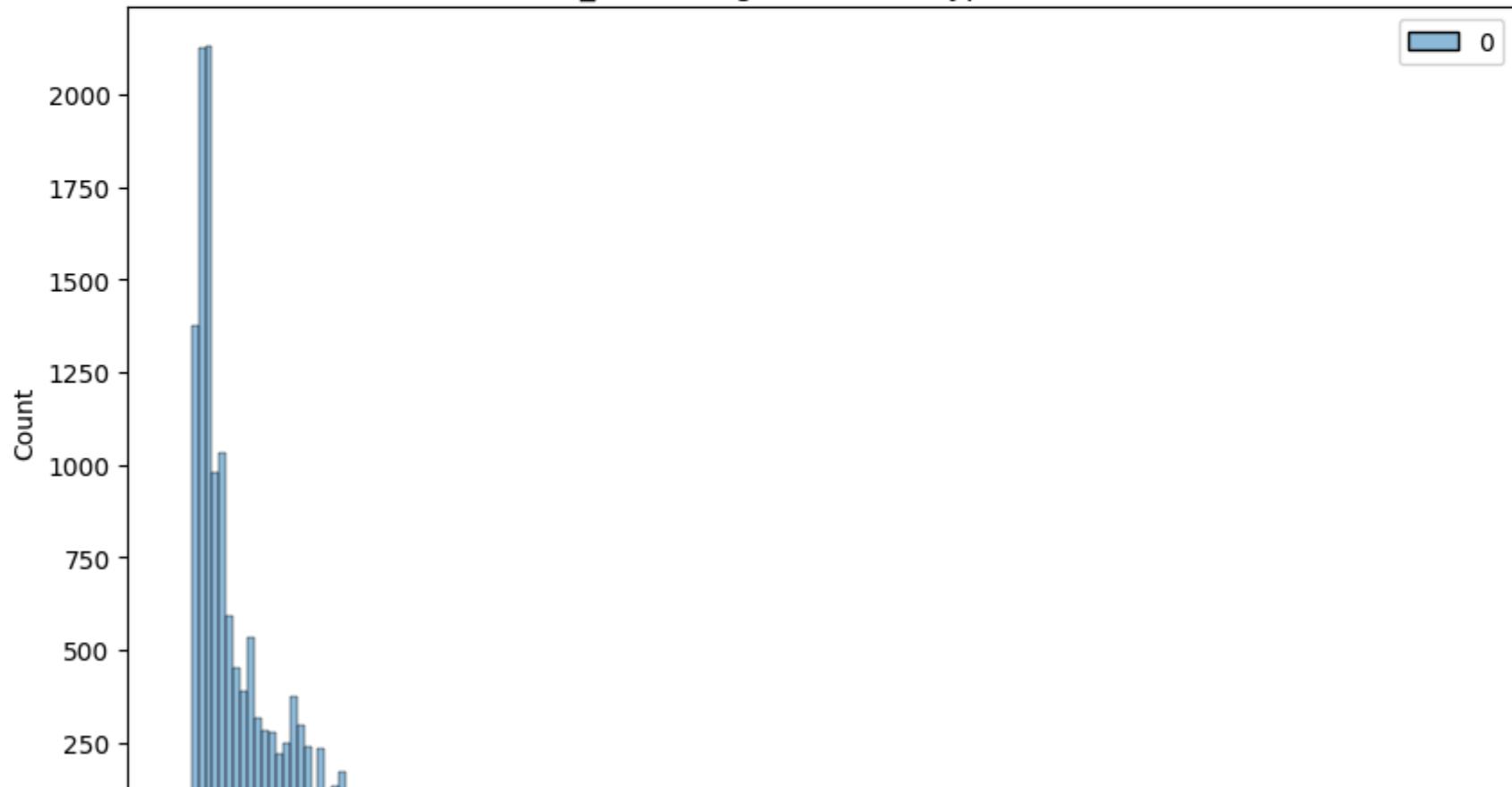
Standardized 0 717.0

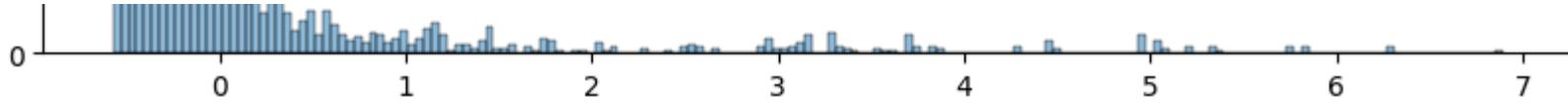
1 68.0
2 1740.0
3 15.0
4 117.0

...

14812 62.0
14813 12.0
14814 48.0
14815 179.0
14816 68.0

Name: osrm_time, Length: 14817, dtype: float32 column





In [240...]

```
plt.figure(figsize = (10, 6))
scaler = StandardScaler()
scaled = scaler.fit_transform(data_2['osrm_distance'].to_numpy().reshape(-1, 1))
sns.histplot(scaled)
plt.title(f"Standardized {data_2['osrm_distance']} column")
plt.plot()
```

Out[240...]

[]

Standardized 0 991.352295

1 85.111000

2 2354.066650

3 19.680000

4 146.791794

...

14812 73.462997

14813 16.088200

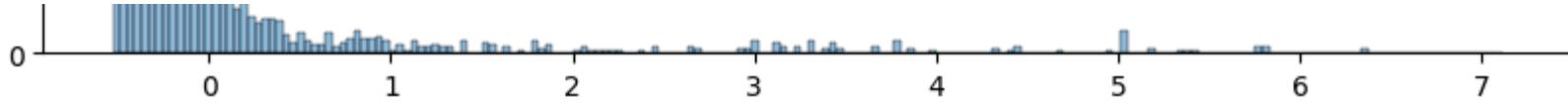
14814 58.903702

14815 171.110306

14816 80.578705

Name: osrm_distance, Length: 14817, dtype: float32 column





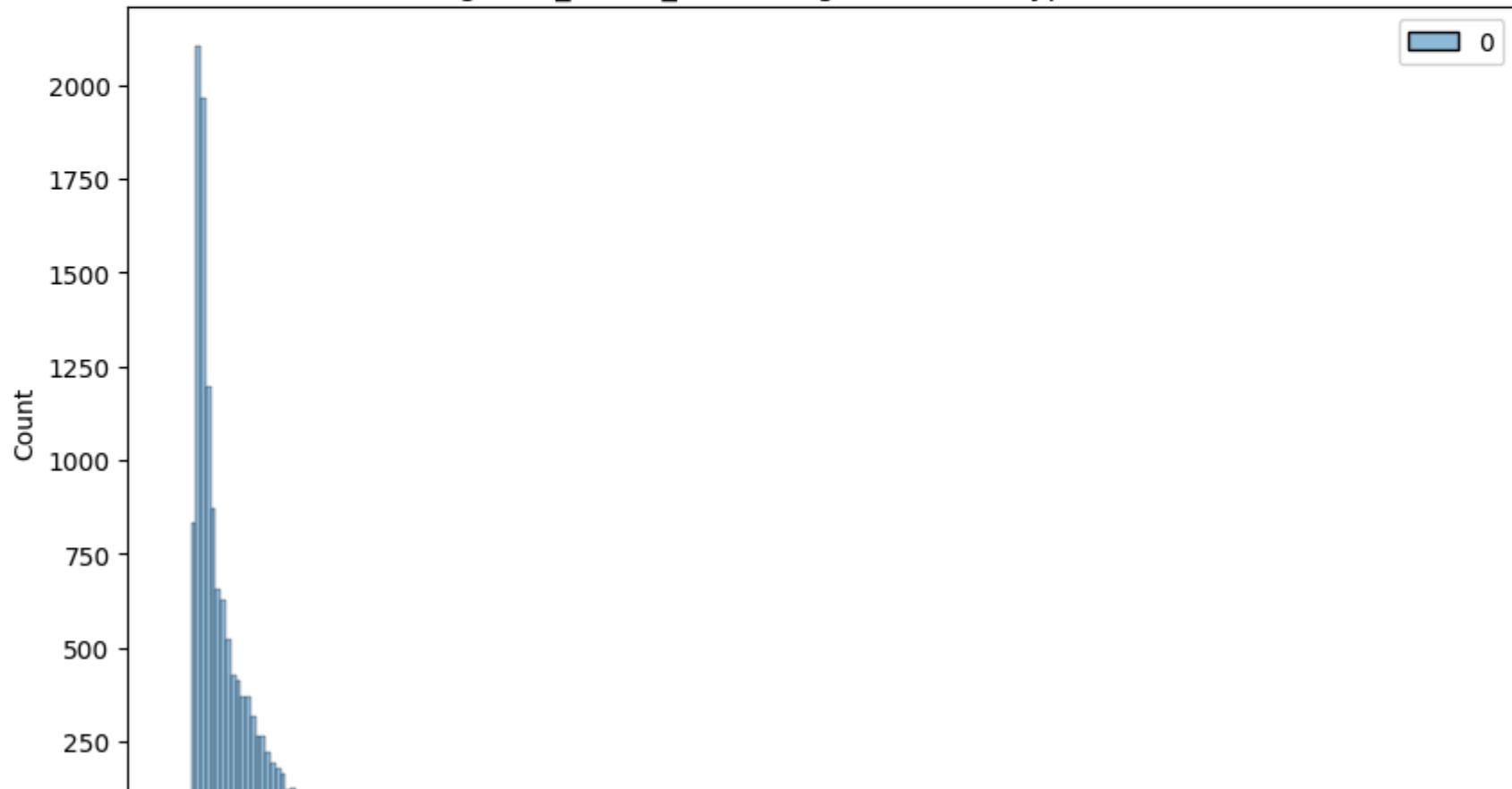
In [241...]

```
plt.figure(figsize = (10, 6))
scaler = StandardScaler()
scaled = scaler.fit_transform(data_2['segment_actual_time'].to_numpy().reshape(-1, 1))
sns.histplot(scaled)
plt.title(f"Standardized {data_2['segment_actual_time']} column")
plt.plot()
```

Out[241...]

[]

```
Standardized 0      1548.0
              1      141.0
              2     3308.0
              3      59.0
              4     340.0
...
              14812    82.0
              14813   21.0
              14814   281.0
              14815   258.0
              14816   274.0
Name: segment_actual_time, Length: 14817, dtype: float32 column
```





In [242...]

```
plt.figure(figsize = (10, 6))
scaler = StandardScaler()
scaled = scaler.fit_transform(data_2['segment_osrm_time'].to_numpy().reshape(-1, 1))
sns.histplot(scaled)
plt.title(f"Standardized {data_2['segment_osrm_time']} column")
plt.plot()
```

Out[242...]

[]

Standardized 0 1008.0

1 65.0

2 1941.0

3 16.0

4 115.0

...

14812 62.0

14813 11.0

14814 88.0

14815 221.0

14816 67.0

Name: segment_osrm_time, Length: 14817, dtype: float32 column





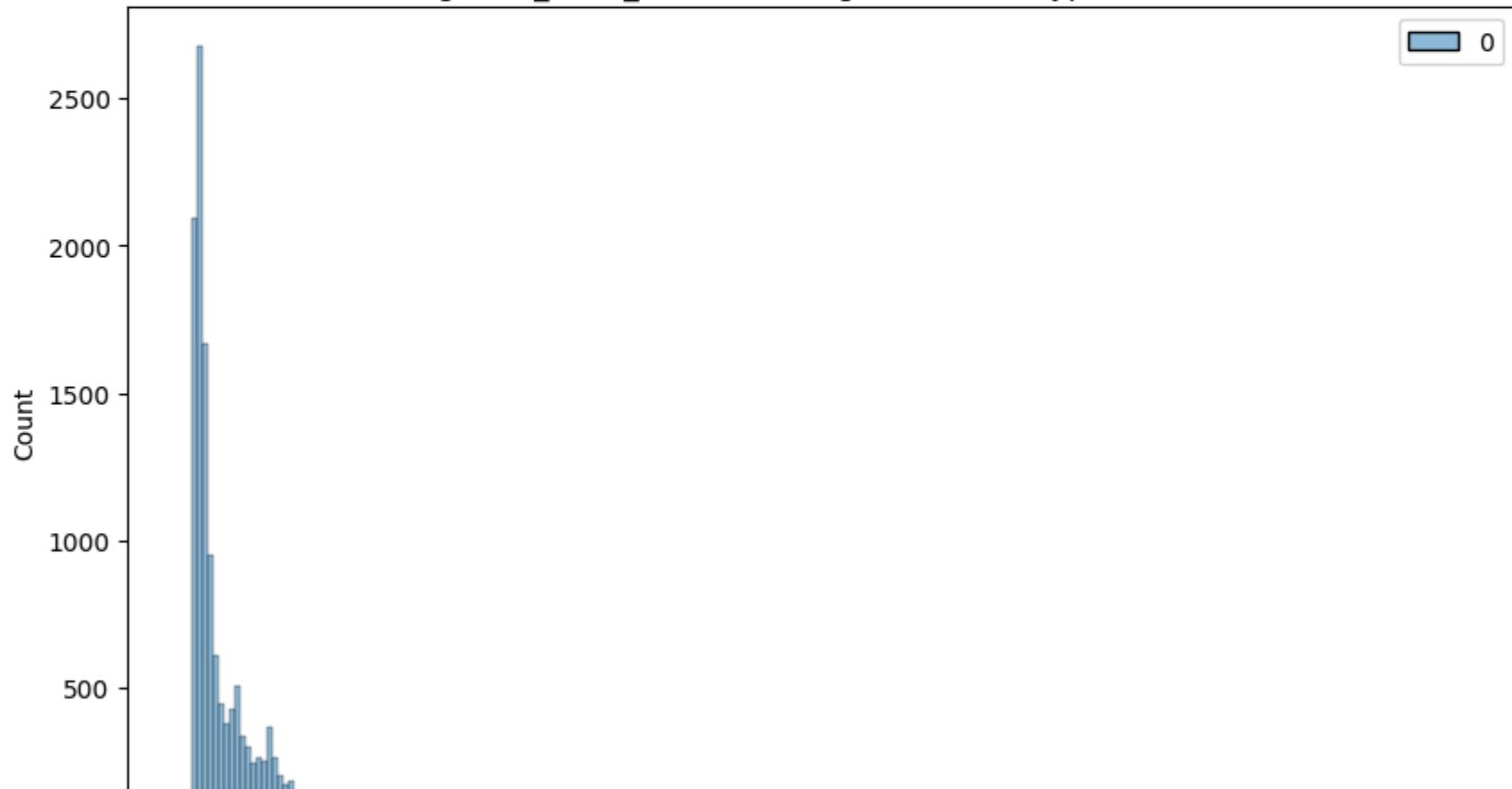
In [243...]

```
plt.figure(figsize = (10, 6))
scaler = StandardScaler()
scaled = scaler.fit_transform(data_2['segment_osrm_distance'].to_numpy().reshape(-1, 1))
sns.histplot(scaled)
plt.title(f"Standardized {data_2['segment_osrm_distance']} column")
plt.plot()
```

Out[243...]

[]

```
Standardized 0      1320.473267
              1      84.189400
              2     2545.267822
              3     19.876600
              4    146.791901
              ...
              14812   64.855103
              14813   16.088299
              14814   104.886597
              14815   223.532394
              14816   80.578705
Name: segment_osrm_distance, Length: 14817, dtype: float32 column
```





Business Insights

The data is given from the period '2018-09-12 00:00:16' to '2018-10-08 03:00:24'.

There are about 14817 unique trip IDs, 1508 unique source centers, 1481 unique destination_centers, 690 unique source cities, 806 unique destination cities.

Most of the data is for testing than for training.

Most common route type is Carting.

The names of 14 unique location ids are missing in the data.

The number of trips start increasing after the noon, becomes maximum at 10 P.M and then start decreasing.

Maximum trips are created in the 38th week.

Most orders come mid-month. That means customers usually make more orders in the mid of the month.

Most orders are sourced from the states like Maharashtra, Karnataka, Haryana, Tamil Nadu, Telangana

Maximum number of trips originated from Mumbai city followed by Gurgaon Delhi, Bengaluru and Bhiwandi. That means that the seller base is strong in these cities.

Maximum number of trips ended in Maharashtra state followed by Karnataka, Haryana, Tamil Nadu and Uttar Pradesh. That means that the number of orders placed in these states is significantly high.

Maximum number of trips ended in Mumbai city followed by Bengaluru, Gurgaon, Delhi and Chennai. That means that the number of orders placed in these cities is significantly high.

Most orders in terms of destination are coming from cities like bengaluru, mumbai, gurgaon, bangalore, Delhi.

Features start_scan_to_end_scan and od_total_time(created feature) are statistically similar.

Features actual_time & osrm_time are statically different.

Features start_scan_to_end_scan and segment_actual_time are statistically similar.

Features osrm_distance and segment_osrm_distance are statistically different from each other.

Both the osrm_time & segment_osrm_time are not statistically same.

Recommendations

The OSRM trip planning system needs to be improved. Discrepancies need to be catered to for transporters, if the routing engine is configured for optimum results.

osrm_time and actual_time are different. Team needs to make sure this difference is reduced, so that better delivery time prediction can be made and it becomes convenient for the customer to expect an accurate delivery time.

The osrm distance and actual distance covered are also not same i.e. maybe the delivery person is not following the predefined route which may lead to late deliveries or the osrm devices is not properly predicting the route based on distance, traffic and other factors. Team needs to look into it.

Most of the orders are coming from/reaching to states like Maharashtra, Karnataka, Haryana and Tamil Nadu. The existing corridors can be further enhanced to improve the penetration in these areas.

Customer profiling of the customers belonging to the states Maharashtra, Karnataka, Haryana, Tamil Nadu and Uttar Pradesh has to be done to get to know why major orders are coming from these states and to improve customers' buying and delivery experience.

From state point of view, we might have very heavy traffic in certain states and bad terrain conditions in certain states. This will be a good indicator to plan and cater to demand during peak festival seasons.

In []:

In []:

