

T **B** **I** <> **alink** **unlink** **bold** **italic** **list-item** **list-item** **list-item** **list-item** **list-item** **list-item** **list-item** **list-item**

```
# <h1 align='center'> <font color='black'><font size=7> 🚚 DELHIVERY  
CaseStudy 🚚 </font> </font></h1>  
<h1 align='center'><font color='black'><font size=6>Feature Engineer  
</font></h1>  
<h2 align='right'>Analysed by : <font color='red'><b> Sayyed Asif Ri  
font></h2>  
  
<kbd>![Delhivery-Success-story-startuptalky1.jpg]  
(attachment:edaa1eb4-32bc-4ce8-a3d4-67ab6a015f9c.jpg)  
![image-asset.gif](attachment:a565dc86-3065-4af7-830f-5044490f3034.g
```



DELHIVERY -

Business CaseStudy 🚚

Feature Engineering

Analysed by : **Sayyed Asif Rizvi**

▼ Introduction:

- Delhivery, established in 2011, is India's foremost logistics and supply chain service provider, offering a comprehensive range of solutions including express parcel transportation, warehousing, and last-mile delivery.
- Leveraging advanced technology and a vast delivery network, Delhivery efficiently manages nationwide movement of goods, earning trust across businesses of all sizes for its dedication to innovation and customer satisfaction.
- As the largest fully integrated player in India by revenue in Fiscal 2021, Delhivery aims to lead the industry by pioneering the commerce operating system, driven by top-tier infrastructure, logistics operations, and innovative data intelligence initiatives led by its Data team.

◆ Why this case study?

Delhivery aims to establish itself as the premier player in the logistics industry. This case study is of paramount importance as it aligns with the company's core objectives and operational excellence.

It provides a practical framework for understanding and processing data, which is integral to their operations. By leveraging data engineering pipelines and data analysis techniques, Delhivery can achieve several critical goals.

First, it allows them to ensure data integrity and quality by addressing missing values and structuring the dataset appropriately.

Second, it enables the extraction of valuable features from raw data, which can be utilized for building accurate forecasting models.

Moreover, it facilitates the identification of patterns, insights, and actionable recommendations crucial for optimizing their logistics operations.

By conducting hypothesis testing and outlier detection, Delhivery can refine their processes and further enhance the quality of service they provide.

How can you help here?

The company wants to understand and process the data coming out of data engineering pipelines:

- Clean, sanitize and manipulate data to get useful features out of raw fields
- Make sense out of the raw data and help the data science team to build forecasting models on it.

▼ Features of the dataset:

- Column Profiling:

| Feature | Description |
|------------------------|---|
| data | tells 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 |
| a. FTL—Full Truck Load | FTL shipments get to the destination sooner, as the truck is making no other pickups or drop-offs along the way |
| b. Carting | Handling system consisting of small vehicles (carts) |

| Feature | Description |
|--------------------------------|---|
| trip_uuid | Unique ID given to a particular trip (A 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 | An open-source routing engine time calculator which computes the shortest path between points in a given map (Includes usual traffic, distance through major and minor roads) (|
| osrm_distance | An open-source routing engine which computes the shortest path between points in a given map (Includes usual traffic, distance through major and minor roads) (|
| factor | Unknown field |
| segment_actual_time | This is a segment time. Time taken by the subset of the package delivery |
| segment_osrm_time | This is the OSRM segment time. Time taken by the subset of the package delivery |
| segment_osrm_distance | This is the OSRM distance. Distance covered by subset of the package delivery |
| segment_factor | Unknown field |

```

1 # importing the required modules and packages
2
3 import pandas as pd
4 import numpy as np
5 import matplotlib.pyplot as plt
6 %matplotlib inline
7 import seaborn as sns
8 import re
9 from scipy.stats import norm,zscore,boxcox,probplot
10 from scipy.stats import ttest_ind,ttest_rel,mannwhitneyu,wilcoxon
11 from scipy.stats import shapiro,levene,kstest,anderson
12 import statsmodels.api as sm
13 from sklearn.impute import SimpleImputer
14 from sklearn.preprocessing import StandardScaler , MinMaxScaler , OneHotEncoder
15 import warnings
16 warnings.filterwarnings('ignore')

```

```

1 # pd_reading the data
2 delhivery_data = pd.read_csv('delhivery_data.csv')

```

```

1 # setting the option of displaying all the columns
2 pd.set_option('display.max_columns', 50)

```

```

1 # making a deep copy for backup
2 dd = delhivery_data.copy()
3 dd.head()

```

| | data | trip_creation_time | route_schedule_uuid | route_type | trip_uuid | source_center | source_name | destination_ |
|---|----------|-------------------------------|---|------------|--------------------|---------------|----------------------------|--------------|
| 0 | training | 2018-09-20 02:35:36.476840 | thanos::sroute:eb7bfc78-b351-4c0e-a951-fa3d5c3... | Carting | 153741093647649320 | IND388121AAA | Anand_VUNagar_DC (Gujarat) | IND3886 |
| 1 | training | 2018-09-20 02:35:36.476840 | thanos::sroute:eb7bfc78-b351-4c0e-a951-fa3d5c3... | Carting | 153741093647649320 | IND388121AAA | Anand_VUNagar_DC (Gujarat) | IND3886 |
| 2 | training | 2018-09-20 02:35:36.476840 | thanos::sroute:eb7bfc78-b351-4c0e-a951-fa3d5c3... | Carting | 153741093647649320 | IND388121AAA | Anand_VUNagar_DC (Gujarat) | IND3886 |
| 3 | training | 2018-09-20 02:35:36.476840 | thanos::sroute:eb7bfc78-b351-4c0e-a951-fa3d5c3... | Carting | 153741093647649320 | IND388121AAA | Anand_VUNagar_DC (Gujarat) | IND3886 |
| 4 | training | 2018-09-20 02:35:36.476840 | thanos::sroute:eb7bfc78-b351-4c0e-a951-fa3d5c3... | Carting | 153741093647649320 | IND388121AAA | Anand_VUNagar_DC (Gujarat) | IND3886 |

Exploration of data :

```
1 dd.shape
```

↳ (144867, 24)

```
1 dd.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')

```
1 dd.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

⌄ Statistical Summary

```
1 dd.describe().T
```

| | 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 |

```
1 dd.describe(include=object).T
```

| | count | unique | | top | freq |
|----------------------------|--------|--------|---|----------|--------|
| data | 144867 | 2 | | training | 104858 |
| trip_creation_time | 144867 | 14817 | 2018-09-28 05:23:15.359220 | 101 | |
| route_schedule_uuid | 144867 | 1504 | thanos::sroute:4029a8a2-6c74-4b7e-a6d8-f9e069f... | 1812 | |
| route_type | 144867 | 2 | | FTL | 99660 |
| trip_uuid | 144867 | 14817 | trip-153811219535896559 | 101 | |
| source_center | 144867 | 1508 | IND000000ACB | 23347 | |
| source_name | 144574 | 1498 | Gurgaon_Bilaspur_HB (Haryana) | 23347 | |
| destination_center | 144867 | 1481 | IND000000ACB | 15192 | |
| destination_name | 144606 | 1468 | Gurgaon_Bilaspur_HB (Haryana) | 15192 | |
| od_start_time | 144867 | 26369 | 2018-09-21 18:37:09.322207 | 81 | |
| od_end_time | 144867 | 26369 | 2018-09-24 09:59:15.691618 | 81 | |
| cutoff_timestamp | 144867 | 93180 | 2018-09-24 05:19:20 | 40 | |

✖️ Duplicate Detection

```
1 dd[dd.duplicated()]
```

| | data | trip_creation_time | route_schedule_uuid | route_type | trip_uuid | source_center | source_name | destination_center | destination_name |
|--|------|--------------------|---------------------|------------|-----------|---------------|-------------|--------------------|------------------|
| | | | | | | | | | |

✖️ Insights

- The dataset does not contain any duplicates.

◆ ? Null Detection

```
1 dd.isna().any()
```

| | |
|---------------------------------------|-------|
| data | False |
| trip_creation_time | False |
| route_schedule_uuid | False |
| route_type | False |
| trip_uuid | False |
| source_center | False |
| source_name | True |
| destination_center | False |
| destination_name | True |
| od_start_time | False |
| od_end_time | False |
| start_scan_to_end_scan | False |
| is_cutoff | False |
| cutoff_factor | False |
| cutoff_timestamp | False |
| actual_distance_to_destination | False |
| actual_time | False |
| osrm_time | False |
| osrm_distance | False |
| factor | False |
| segment_actual_time | False |
| segment_osrm_time | False |
| segment_osrm_distance | False |
| segment_factor | False |
| dtype: bool | |

```
1 dd.isnull().sum()
```

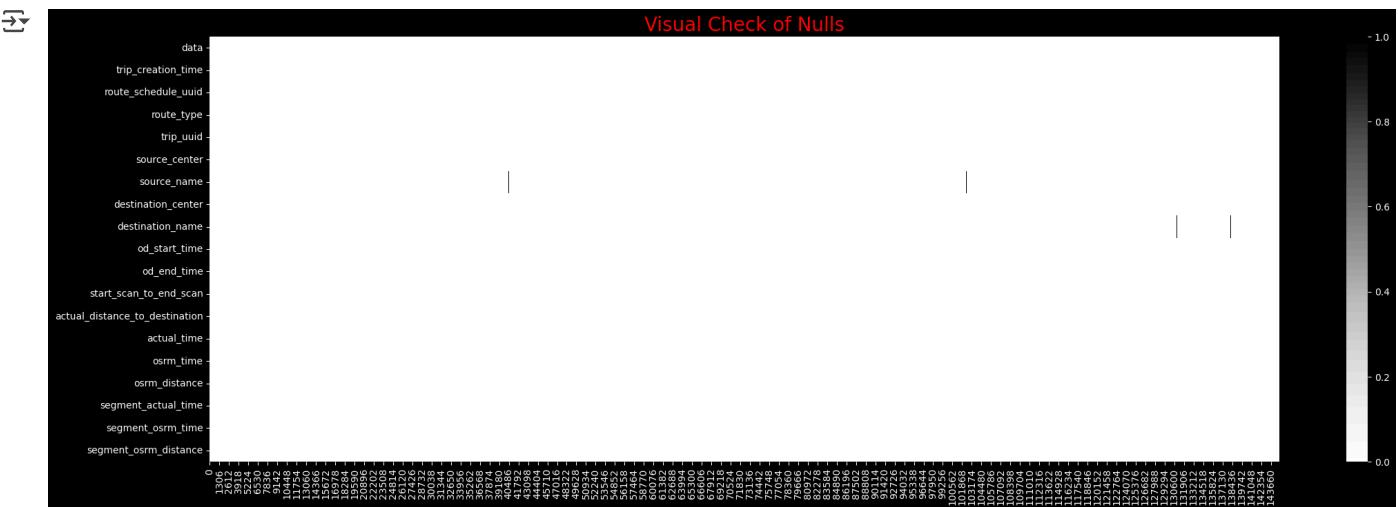
| | |
|-------------------------------|-----|
| data | 0 |
| trip_creation_time | 0 |
| route_schedule_uuid | 0 |
| route_type | 0 |
| trip_uuid | 0 |
| source_center | 0 |
| source_name | 293 |
| destination_center | 0 |
| destination_name | 261 |
| od_start_time | 0 |
| od_end_time | 0 |
| start_scan_to_end_scan | 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
```

```
1 def missing_data(df):
2     total_missing_df = df.isnull().sum().sort_values(ascending = False)
3     percent_missing_df = (df.isnull().sum()/df.isna().count()*100).sort_values(ascending=False) # -----> /len(dd)
4     missing_data_df = pd.concat([total_missing_df, percent_missing_df], axis=1, keys=['Total', 'Percent'])
5     return missing_data_df
6
7 missing_pct = missing_data(dd)
8 missing_pct[missing_pct['Total']>0]
```

| | Total | Percent |
|------------------|-------|----------|
| source_name | 293 | 0.202254 |
| destination_name | 261 | 0.180165 |

```
1 plt.figure(figsize=(25,8))
2 plt.style.use('dark_background')
3 sns.heatmap(dd.isnull().T,cmap='Greys')
4 plt.title('Visual Check of Nulls', fontsize=20,color='r')
5 plt.show()
```



```
1 # Dropping unknown fields
2
3 unknown_fields = ['is_cutoff', 'cutoff_factor', 'cutoff_timestamp', 'factor', 'segment_factor']
4 dd = dd.drop(columns = unknown_fields)

1 dd.sample()
```

| | data | trip_creation_time | route_schedule_uuid | route_type | trip_uuid | source_center | source_name | destinat |
|--------|------|----------------------------|--|------------|--------------------|-------------------|-------------------------------|----------|
| 133488 | test | 2018-09-30 05:56:48.299467 | thanos::sroute:6be6529bf2ad-4714-b7ab-ac58f24... | FTL | 153828700829921150 | trip-IND000000ACB | Gurgaon_Bilaspur_HB (Haryana) | INC |

```
1 dd.shape
```

```
2 (144867, 19)
```

```

1 #checking the unique values for columns
2 for _ in dd.columns:
3     print()
4     print(f'Total Unique Values in {_} column are :- {dd[_].nunique()}')
5     print(f'Unique Values in {_} column are :-\n {dd[_].unique()}')
6     print()
7     print('-'*120)
→ 6.310e+02 6.780e+02 4.010e+02 5.970e+02 9.270e+02 1.167e+03
  1.115e+03 3.260e+02 7.150e+02 1.930e+02 1.395e+03 9.670e+02
  5.370e+02 6.040e+02 3.210e+02 3.770e+02 8.510e+02 1.016e+03
  1.031e+03 4.550e+02 6.220e+02 3.890e+02 5.830e+02 -4.800e+01
  4.180e+02 1.047e+03 6.930e+02 9.700e+02 7.040e+02 7.850e+02
-2.000e+00 2.491e+03 9.460e+02 4.650e+02 2.541e+03 1.122e+03
  3.051e+03 9.740e+02 9.780e+02 9.040e+02 -1.600e+01 8.530e+02
  4.790e+02 1.148e+03 5.720e+02 4.250e+02 5.530e+02 4.060e+02
-7.000e+00 3.070e+02 1.093e+03 4.630e+02 9.390e+02 5.450e+02
  1.325e+03 9.150e+02 5.460e+02 7.530e+02 5.290e+02 4.370e+02
  5.200e+02 8.300e+02 1.677e+03 1.020e+03 7.480e+02 4.880e+02
  6.130e+02 9.510e+02 3.740e+02 7.360e+02 9.330e+02 5.790e+02
  2.625e+03 7.520e+02 -1.500e+01 1.192e+03 6.640e+02 1.320e+03
  2.870e+02 3.700e+02 1.104e+03]

```

Total Unique Values in segment_osrm_time column are :- 214

Unique Values in segment_osrm_time column are :-

```

[1.100e+01 9.000e+00 7.000e+00 1.200e+01 5.000e+00 6.000e+00 1.000e+01
 2.400e+01 2.700e+01 2.600e+01 1.400e+01 1.500e+01 3.000e+01 1.800e+01
 3.800e+01 3.700e+01 2.500e+01 1.700e+01 2.200e+01 3.600e+01 3.200e+01
 1.600e+01 7.000e+01 3.500e+01 4.500e+01 1.300e+01 0.000e+00 8.000e+00
 2.000e+00 1.900e+01 2.300e+01 2.800e+01 2.000e+01 2.100e+01 3.300e+01
 3.400e+01 8.100e+01 3.000e+00 4.400e+01 1.000e+00 4.000e+00 3.900e+01
 4.000e+01 2.900e+01 5.300e+01 3.100e+01 7.500e+01 7.900e+01 9.700e+01
 4.800e+01 4.100e+01 4.300e+01 5.000e+01 5.400e+01 6.800e+01 4.200e+01
 4.600e+01 6.000e+01 5.800e+01 7.600e+01 4.900e+01 1.300e+02 7.800e+01
 6.700e+01 6.400e+01 5.500e+01 5.100e+01 1.420e+02 7.700e+01 7.100e+01
 5.600e+01 6.600e+01 5.900e+01 9.200e+01 6.200e+01 5.200e+01 5.700e+01
 8.000e+01 4.700e+01 7.400e+01 6.500e+01 8.800e+01 7.200e+01 6.900e+01
 7.300e+01 8.400e+01 8.200e+01 8.300e+01 1.540e+02 9.100e+01 6.100e+01
 9.400e+01 1.220e+02 6.300e+01 2.180e+02 9.800e+01 8.700e+01 3.830e+02
 9.000e+01 1.080e+02 9.300e+01 8.600e+01 8.900e+01 1.020e+02 1.600e+02
 2.210e+02 1.050e+02 1.330e+02 1.000e+02 1.060e+02 4.070e+02 8.500e+01
 1.340e+02 4.690e+02 1.800e+02 2.340e+02 1.490e+02 1.010e+02 1.450e+02
 1.140e+02 1.840e+02 2.270e+02 1.740e+02 1.320e+02 9.900e+01 9.600e+01
 1.310e+02 1.110e+02 1.040e+02 1.750e+02 2.300e+02 9.500e+01 1.250e+02
 2.950e+02 1.560e+02 1.160e+02 1.460e+02 1.410e+02 1.030e+02 1.170e+02
 2.310e+02 2.540e+02 2.200e+02 2.330e+02 1.810e+02 1.210e+02 1.270e+02
 3.700e+02 3.750e+02 1.500e+02 1.070e+02 1.610e+02 2.320e+02 1.090e+02
 1.200e+02 1.100e+02 9.970e+02 1.790e+02 1.130e+02 1.660e+02 9.960e+02
 1.240e+02 2.150e+02 1.570e+02 3.620e+02 1.430e+02 1.150e+02 1.280e+02
 1.700e+02 1.440e+02 2.350e+02 1.510e+02 3.560e+02 1.180e+02 1.390e+02
 1.710e+02 1.290e+02 1.190e+02 1.690e+02 1.630e+02 2.040e+02 1.480e+02
 1.830e+02 4.810e+02 3.410e+02 3.280e+02 2.130e+02 1.890e+02 1.910e+02
 1.400e+02 1.470e+02 2.080e+02 2.860e+02 2.160e+02 1.720e+02 1.380e+02
 1.670e+02 2.940e+02 1.230e+02 1.260e+02 2.110e+02 1.611e+03 2.190e+02
 2.490e+02 1.850e+02 1.580e+02 3.240e+02 1.770e+02 4.530e+02 1.520e+02
 1.760e+02 7.370e+02 1.730e+02 1.032e+03]

```

Total Unique Values in segment_osrm_distance column are :- 113799

Unique Values in segment_osrm_distance column are :-

```
[11.9653 9.759 10.8152 ... 20.7053 18.8885 8.8088]
```

Changing the Datatype of Columns

```
1 dd.sample()
```

| | data | trip_creation_time | route_schedule_uuid | route_type | trip_uuid | source_center | source_name | destina |
|--------|----------|-------------------------------|--|------------|-----------------------------|---------------|--------------------------------|---------|
| 124444 | training | 2018-09-23 02:44:13.665024 | thanos::route:f01c8bbd- 655d-42ea-9abf- 60d5040... | FTL | trip- 153767065366477819 | IND821115AAB | Sasaram_Central_I_2 (Bihar) | IN |

```
1 dd.dtypes
```

| | |
|---------------------|--------|
| → data | object |
| trip_creation_time | object |
| route_schedule_uuid | object |

```

route_type          object
trip_uuid          object
source_center       object
source_name         object
destination_center object
destination_name   object
od_start_time      object
od_end_time        object
start_scan_to_end_scan float64
actual_distance_to_destination float64
actual_time         float64
osrm_time           float64
osrm_distance       float64
segment_actual_time float64
segment_osrm_time  float64
segment_osrm_distance float64
dtype: object

```

```

1 # Converting the datatypes to category for columns like data and route_type as they only have 2 values.
2 dd['data'] = dd['data'].astype('category')
3 dd['route_type'] = dd['route_type'].astype('category')
4
5 # Converting time columns to datetime format
6 datetime_cols = ['trip_creation_time', 'od_start_time', 'od_end_time']
7 for _ in datetime_cols:
8     dd[_] = pd.to_datetime(dd[_])

```

```
1 dd.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   float64
 13  actual_time         144867 non-null   float64
 14  osrm_time           144867 non-null   float64
 15  osrm_distance       144867 non-null   float64
 16  segment_actual_time 144867 non-null   float64
 17  segment_osrm_time  144867 non-null   float64
 18  segment_osrm_distance 144867 non-null   float64
dtypes: category(2), datetime64[ns](3), float64(8), object(6)
memory usage: 19.1+ MB

```

```

float_cols = []
for _ in dd.columns:
    if isinstance(dd[_], 'float64'):
        float_cols.append(_)
float_cols

```

✓ see y it didnt work

```

1 float_cols = []
2 for _ in dd.columns:
3     if dd[_].dtype=='float64':
4         float_cols.append(_)
5 float_cols

```

```

→ ['start_scan_to_end_scan',
  'actual_distance_to_destination',
  'actual_time',
  'osrm_time',
  'osrm_distance',
  'segment_actual_time',
  'segment_osrm_time',
  'segment_osrm_distance']

```

```

1 # reducing the float64 to float32 to save memory
2 for _ in float_cols:
3     dd[_] = dd[_].astype('float32')

1 dd.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   float32
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(8), object(6)
memory usage: 14.6+ MB

```

↳ 🌐 Insights:

- Earlier the dataset was using 25.6+ MB of memory but now it has been reduced to 14.6 + MB. Around 40.63 % reduction in the memory usage.

```

1 # Time period of data
2 dd['trip_creation_time'].max(), dd['trip_creation_time'].min(), dd['trip_creation_time'].max()-dd['trip_creation_time'].min()

→ (Timestamp('2018-10-03 23:59:42.701692'),
 Timestamp('2018-09-12 00:00:16.535741'),
 Timedelta('21 days 23:59:26.165951'))

1 # Time period of data
2 dd['od_start_time'].max(), dd['od_start_time'].min(), dd['od_start_time'].max() - dd['od_start_time'].min()

→ (Timestamp('2018-10-06 04:27:23.392375'),
 Timestamp('2018-09-12 00:00:16.535741'),
 Timedelta('24 days 04:27:06.856634'))

1 # Time period of data
2 dd['od_end_time'].max(), dd['od_end_time'].min(), dd['od_end_time'].max() - dd['od_end_time'].min()

→ (Timestamp('2018-10-08 03:00:24.353479'),
 Timestamp('2018-09-12 00:50:10.814399'),
 Timedelta('26 days 02:10:13.539080'))

1 data_time_frame = dd['od_end_time'].max() - dd['trip_creation_time'].min()
2 data_time_frame

→ Timedelta('26 days 03:00:07.817738')

```

↳ 🔴 Null Treatment:

- Replace null values in 'source_name' and 'destination_name' columns with 'unknown' through scikit imputation

```

columns_to_impute = ['source_name', 'destination_name']
imputer = SimpleImputer(strategy='constant', fill_value='unknown')
dd[columns_to_impute] = imputer.fit_transform(dd[columns_to_impute])

but 'unknown' transactions will be more and to be omitted while analyzing ... Hence no use of imputing....

```

```
1 dd[(dd.source_name.isna())&(dd.destination_name.isna())]
```

| | data | trip_creation_time | route_schedule_uuid | route_type | | trip_uuid | source_center | source_name | destination_cer |
|-------|----------|-------------------------------|---|------------|-----------------------------|--------------|---------------|-------------|-----------------|
| 68006 | training | 2018-09-26 22:21:56.619259 | thanos::sroute:cfb575b8- df26-48f5-8427- 6f48f9d... | FTL | trip- 153800051661903546 | IND331022A1B | NaN | IND331001 | |
| 68007 | training | 2018-09-26 22:21:56.619259 | thanos::sroute:cfb575b8- df26-48f5-8427- 6f48f9d... | FTL | trip- 153800051661903546 | IND331022A1B | NaN | IND331001 | |
| 68008 | training | 2018-09-26 22:21:56.619259 | thanos::sroute:cfb575b8- df26-48f5-8427- 6f48f9d... | FTL | trip- 153800051661903546 | IND331022A1B | NaN | IND331001 | |

```
1 dd[dd.source_name.isna()]
```

| | data | trip_creation_time | route_schedule_uuid | route_type | | trip_uuid | source_center | source_name | destination_c |
|--------|----------|-------------------------------|--|------------|-----------------------------|--------------|---------------|-------------|---------------|
| 112 | training | 2018-09-25 08:53:04.377810 | thanos::sroute:4460a38d- ab9b-484e-bd4e- f4201d0... | FTL | trip- 153786558437756691 | IND342902A1B | NaN | IND30201 | |
| 113 | training | 2018-09-25 08:53:04.377810 | thanos::sroute:4460a38d- ab9b-484e-bd4e- f4201d0... | FTL | trip- 153786558437756691 | IND342902A1B | NaN | IND30201 | |
| 114 | training | 2018-09-25 08:53:04.377810 | thanos::sroute:4460a38d- ab9b-484e-bd4e- f4201d0... | FTL | trip- 153786558437756691 | IND342902A1B | NaN | IND30201 | |
| 115 | training | 2018-09-25 08:53:04.377810 | thanos::sroute:4460a38d- ab9b-484e-bd4e- f4201d0... | FTL | trip- 153786558437756691 | IND342902A1B | NaN | IND30201 | |
| 116 | training | 2018-09-25 08:53:04.377810 | thanos::sroute:4460a38d- ab9b-484e-bd4e- f4201d0... | FTL | trip- 153786558437756691 | IND342902A1B | NaN | IND30201 | |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| 144484 | test | 2018-10-03 09:06:06.690094 | thanos::sroute:cbeff3b6a- 79ea-4d5e-a215- b558a70... | FTL | trip- 153855756668984584 | IND282002AAD | NaN | IND47400 | |
| 144485 | test | 2018-10-03 09:06:06.690094 | thanos::sroute:cbeff3b6a- 79ea-4d5e-a215- b558a70... | FTL | trip- 153855756668984584 | IND282002AAD | NaN | IND47400 | |
| 144486 | test | 2018-10-03 09:06:06.690094 | thanos::sroute:cbeff3b6a- 79ea-4d5e-a215- b558a70... | FTL | trip- 153855756668984584 | IND282002AAD | NaN | IND47400 | |
| 144487 | test | 2018-10-03 09:06:06.690094 | thanos::sroute:cbeff3b6a- 79ea-4d5e-a215- b558a70... | FTL | trip- 153855756668984584 | IND282002AAD | NaN | IND47400 | |
| 144488 | test | 2018-10-03 09:06:06.690094 | thanos::sroute:cbeff3b6a- 79ea-4d5e-a215- b558a70... | FTL | trip- 153855756668984584 | IND282002AAD | NaN | IND47400 | |

293 rows × 19 columns

```
1 dd[dd.destination_name.isna()]
```

| | data | trip_creation_time | route_schedule_uuid | route_type | trip_uuid | source_center | source_name | desti |
|--------|----------|-------------------------------|--|------------|--------------------|---------------|------------------------------------|-------|
| 110 | training | 2018-09-25 08:53:04.377810 | thanos::sroute:4460a38d-ab9b-484e-bd4e-f4201d0... | FTL | 153786558437756691 | IND342601AAA | Piparcity_BsstdDPP_D (Rajasthan) | |
| 111 | training | 2018-09-25 08:53:04.377810 | thanos::sroute:4460a38d-ab9b-484e-bd4e-f4201d0... | FTL | 153786558437756691 | IND342601AAA | Piparcity_BsstdDPP_D (Rajasthan) | |
| 982 | test | 2018-10-01 20:56:18.155260 | thanos::sroute:d0ebdacd-e09b-47d3-be77-c9c4a05... | FTL | 153842737815495661 | IND573103AAA | Arsikere_HsnRdDPP_D (Karnataka) | |
| 983 | test | 2018-10-01 20:56:18.155260 | thanos::sroute:d0ebdacd-e09b-47d3-be77-c9c4a05... | FTL | 153842737815495661 | IND573103AAA | Arsikere_HsnRdDPP_D (Karnataka) | |
| 4882 | training | 2018-09-24 07:18:06.087341 | thanos::sroute:2f43f11e-d3ba-4590-9355-82928e1... | FTL | 153777348608709328 | IND202001AAB | Aligarh_KhirByps_I (Uttar Pradesh) | |
| ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 144478 | test | 2018-10-03 09:06:06.690094 | thanos::sroute:cbeef3b6a-79ea-4d5e-a215-b558a70... | FTL | 153855756668984584 | IND000000ACB | Gurgaon_Bilaspur_HB (Haryana) | |
| 144479 | test | 2018-10-03 09:06:06.690094 | thanos::sroute:cbeef3b6a-79ea-4d5e-a215-b558a70... | FTL | 153855756668984584 | IND000000ACB | Gurgaon_Bilaspur_HB (Haryana) | |
| 144480 | test | 2018-10-03 09:06:06.690094 | thanos::sroute:cbeef3b6a-79ea-4d5e-a215-b558a70... | FTL | 153855756668984584 | IND000000ACB | Gurgaon_Bilaspur_HB (Haryana) | |
| 144481 | test | 2018-10-03 09:06:06.690094 | thanos::sroute:cbeef3b6a-79ea-4d5e-a215-b558a70... | FTL | 153855756668984584 | IND000000ACB | Gurgaon_Bilaspur_HB (Haryana) | |
| 144482 | test | 2018-10-03 09:06:06.690094 | thanos::sroute:cbeef3b6a-79ea-4d5e-a215-b558a70... | FTL | 153855756668984584 | IND000000ACB | Gurgaon_Bilaspur_HB (Haryana) | |

261 rows × 19 columns

```
1 ddd = dd.copy()
```

```
1 missing_source_name = ddd.loc[ddd['source_name'].isnull(), 'source_center']
2 missing_source_name
```



```
→ array(['IND342902A1B', 'IND577116AAA', 'IND282002AAD', 'IND465333A1B',
       'IND841301AAC', 'IND509103AAC', 'IND126116AAA', 'IND331022A1B',
       'IND507301AAC', 'IND509103AAC', 'IND126116AAA', 'IND331022A1B'])
```

```
1 missing_destination_name = ddd.loc[ddd['destination_name'].isnull(), 'destination_center'].unique()
2 missing_destination_name
```

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

```
1 # checking if element of one np.array isin another np.array
2
3 #np.all(df.loc[df['source_name'].isnull(), 'source_center'].isin(missing_destination_name))
4
5 np.in1d(missing_source_name, missing_destination_name).all()
```

```
False

1 for _ in missing_source_name:
2     unique_source_name = ddd.loc[ddd['source_center'] == _, 'source_name'].unique()
3     if pd.isna(unique_source_name):
4         print("Source Center :", _, " - " * 5, "Source Name :", 'NA')
5     else :
6         print("Source Center :", . . . " - " * 5, "Source Name :", unique_source_name)
```

→ Source Center : IND342902A1B ----- Source Name : NA
Source Center : IND577116AAA ----- Source Name : NA
Source Center : IND282002AAD ----- Source Name : NA
Source Center : IND465333A1B ----- Source Name : NA
Source Center : IND841301AAC ----- Source Name : NA
Source Center : IND509103AAC ----- Source Name : NA

```
Source Center : IND126116AAA ----- Source Name : NA
Source Center : IND331022A1B ----- Source Name : NA
Source Center : IND505326AAB ----- Source Name : NA
Source Center : IND852118A1B ----- Source Name : NA
```

```
1 for _ in missing_destination_name:
2     unique_destination_name = ddd.loc[ddd['destination_center'] == _, 'destination_name'].unique()
3     if pd.isna(unique_destination_name):
4         print("Destination Center :", _, " -" * 5, "Destination Name :", 'NA')
5     else :
6         print("Destination Center :", _, " -" * 5, "Destination Name :", unique_destination_name)
```

```
→ Destination Center : IND342902A1B ----- Destination Name : NA
Destination Center : IND577116AAA ----- Destination Name : NA
Destination Center : IND282002AAD ----- Destination Name : NA
Destination Center : IND465333A1B ----- Destination Name : NA
Destination Center : IND841301AAC ----- Destination Name : NA
Destination Center : IND505326AAB ----- Destination Name : NA
Destination Center : IND852118A1B ----- Destination Name : NA
Destination Center : IND126116AAA ----- Destination Name : NA
Destination Center : IND509103AAC ----- Destination Name : NA
Destination Center : IND221005A1A ----- Destination Name : NA
Destination Center : IND250002AAC ----- Destination Name : NA
Destination Center : IND331001A1C ----- Destination Name : NA
Destination Center : IND122015AAC ----- Destination Name : NA
```

```
1 count = 1
2 for i in missing_destination_name:
3     ddd.loc[ddd['destination_center'] == i, 'destination_name'] = ddd.loc[ddd['destination_center'] == i,
4                                         'destination_name'].replace(np.nan, f'location_{count}')
5     count += 1
```

```
1 d = {}
2 for i in missing_source_name:
3     d[i] = ddd.loc[ddd['destination_center'] == i, 'destination_name'].unique()
4 for idx, val in d.items():
5     if len(val) == 0:
6         d[idx] = [f'location_{count}']
7         count += 1
8 d2 = {}
9 for idx, val in d.items():
10    d2[idx] = val[0]
11 for i, v in d2.items():
12     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
```

```
1 for i in missing_source_name:
2     ddd.loc[ddd['source_center'] == i, 'source_name'] = ddd.loc[ddd['source_center'] == i, 'source_name'].replace(np.nan, d2[i])
```

```
1 ddd.source_name.value_counts()
```

```
→ source_name
Gurgaon_Bilaspur_HB (Haryana)      23347
Bangalore_Nelmgla_H (Karnataka)    9975
Bhiwandi_Mankoli_H (Maharashtra)   9088
Pune_Tathawde_H (Maharashtra)       4061
Hyderabad_Shamsabd_H (Telangana)   3340
...
Badkulla_Central_DPP_1 (West Bengal) 1
Kasganj_BnkrGate_D (Uttar Pradesh)  1
Shahjhnpur_NavdaCln_D (Uttar Pradesh) 1
Jaunpur_Katghara_D (Uttar Pradesh)  1
Krishnanagar_AnadiDPP_D (West Bengal) 1
Name: count, Length: 1508, dtype: int64
```

```
1 ddd.destination_name.value_counts()
```

```
→ destination_name
Gurgaon_Bilaspur_HB (Haryana)      15192
Bangalore_Nelmgla_H (Karnataka)    11019
Bhiwandi_Mankoli_H (Maharashtra)   5492
Hyderabad_Shamsabd_H (Telangana)   5142
```

```
Kolkata_Dankuni_HB (West Bengal)      4892
Vijayawada (Andhra Pradesh)          1
Ranaghat_ArckDPP_D (West Bengal)    1
Mumbai_Sanpada_CP (Maharashtra)     1
Delhi_Lajwanti (Delhi)              1
Luxettipet_ShivaDPP_D (Telangana)   1
Name: count, Length: 1481, dtype: int64
```

even if we replace these nulls with some values, those are not gonna have any impact on the data.... so we can drop it as well..

```
1 261+290
```

```
⤵ 551
```

```
1 len(delhivery_data)
```

```
⤵ 144867
```

```
1 144867 - 551
```

```
⤵ 144316
```

```
1 df = dd.dropna()
```

```
1 df.isna().sum().any()
```

```
⤵ False
```

```
1 df.info()
```

```
⤵ <class 'pandas.core.frame.DataFrame'>
Index: 144316 entries, 0 to 144866
Data columns (total 19 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   data             144316 non-null   category
 1   trip_creation_time 144316 non-null   datetime64[ns]
 2   route_schedule_uuid 144316 non-null   object  
 3   route_type        144316 non-null   category
 4   trip_uuid         144316 non-null   object  
 5   source_center     144316 non-null   object  
 6   source_name       144316 non-null   object  
 7   destination_center 144316 non-null   object  
 8   destination_name  144316 non-null   object  
 9   od_start_time     144316 non-null   datetime64[ns]
 10  od_end_time       144316 non-null   datetime64[ns]
 11  start_scan_to_end_scan 144316 non-null   float32
 12  actual_distance_to_destination 144316 non-null   float32
 13  actual_time        144316 non-null   float32
 14  osrm_time          144316 non-null   float32
 15  osrm_distance      144316 non-null   float32
 16  segment_actual_time 144316 non-null   float32
 17  segment_osrm_time  144316 non-null   float32
 18  segment_osrm_distance 144316 non-null   float32
dtypes: category(2), datetime64[ns](3), float32(8), object(6)
memory usage: 15.7+ MB
```

⌄ ⚡ Insights:

- Only two fields have a tiny fraction of missing values, less than 0.05% of the whole dataset.
- Since we have plenty of data to work with, we're choosing to just get rid of the missing values instead of trying to guess them using methods like using the average or most common value.
- I'm dropping the missing values to keep things simple and not mess up how the features are spread out. But if a lot more data was missing, we could have used other methods like guessing based on what's there or using the most common values.

🧙 Exploratory Data Analysis

```
1 cp = ['gray','red','dimgrey','tomato','dimgray','orangered','k','salmon','gray','red','dimgrey','tomato','dimgray','orangered','k',':
```

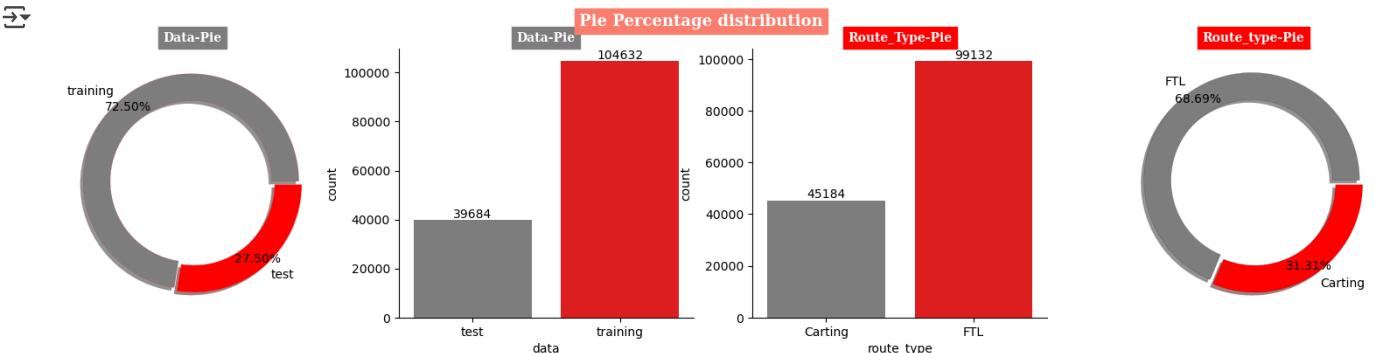
```
1 df.sample()
```

| | data | trip_creation_time | route_schedule_uuid | route_type | trip_uuid | source_center | source_name | destination |
|-------|----------|-------------------------------|---|------------|-------------------------|---------------|------------------------------------|-------------|
| 92389 | training | 2018-09-14 23:36:50.771430 | thanos::sroute:3f001d56-e933-4ba1-a7cf-26828f7... | Carting | trip-153696821077115152 | IND515201AAA | Hindupur_Parigi_D (Andhra Pradesh) | IND51 |

```

1 plt.figure(figsize=(20,4))
2 plt.suptitle('Pie Percentage distribution', fontsize=13, fontfamily='serif', fontweight='bold', backgroundcolor=cp[-1], color='w')
3
4 plt.subplot(141)
5 plt.pie(df['data'].value_counts(), labels=df['data'].value_counts().index, colors=cp, counterclock=True, explode=(0.02, 0.02), autopct=6
6         textprops={'color': 'k', 'fontsize': 10}, shadow=True, radius=1, wedgeprops=dict(edgecolor='r', linewidth=0.1, width=0.25))
7 plt.title('Data-Pie', fontsize=10, fontfamily='serif', fontweight='bold', backgroundcolor=cp[0], color='w')
8
9 plt.subplot(142)
10 a = sns.barplot(x=df['data'].value_counts().index, y=df['data'].value_counts(), palette=cp)
11 a.bar_label(a.containers[0], label_type='edge', fmt='%d')
12 plt.title('Data-Pie', fontsize=10, fontfamily='serif', fontweight='bold', backgroundcolor=cp[0], color='w')
13
14 plt.subplot(143)
15 b = sns.barplot(x=df['route_type'].value_counts().index, y=df['route_type'].value_counts(), palette=cp)
16 b.bar_label(b.containers[0], label_type='edge', fmt='%d')
17 plt.title('Route_Type-Pie', fontsize=10, fontfamily='serif', fontweight='bold', backgroundcolor=cp[1], color='w')
18
19 plt.subplot(144)
20 plt.pie(df['route_type'].value_counts(), labels=df['route_type'].value_counts().index, colors=cp, counterclock=True, explode=(0.03, 0
21         textprops={'color': 'k', 'fontsize': 10}, shadow=True, radius=1, wedgeprops=dict(edgecolor='k', linewidth=0.1, width=0.25))
22 plt.title('Route_type-Pie', fontsize=10, fontfamily='serif', fontweight='bold', backgroundcolor=cp[1], color='w')
23
24 sns.despine()
25 plt.show()

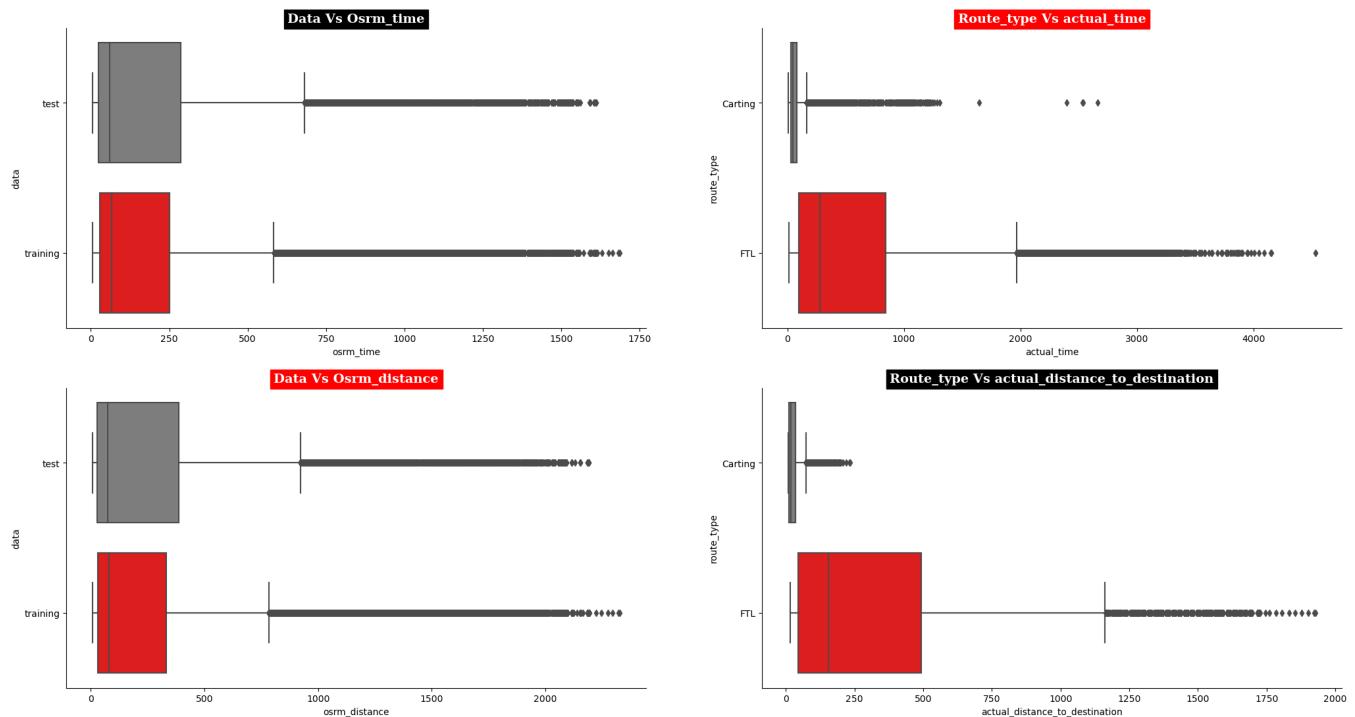
```



```

1 plt.figure(figsize=(25,13))
2 plt.style.use('default')
3 plt.style.use('seaborn-bright')
4
5 plt.subplot(221)
6 sns.boxplot(data=df, y='data', x='osrm_time', palette=cp)
7 plt.title('Data Vs Osrm_time', fontsize=14, fontfamily='serif', fontweight='bold', backgroundcolor='k', color='w')
8
9 plt.subplot(222)
10 sns.boxplot(data=df, y='route_type', x='actual_time', palette=cp)
11 plt.title('Route_type Vs actual_time', fontsize=14, fontfamily='serif', fontweight='bold', backgroundcolor='r', color='w')
12
13 plt.subplot(223)
14 sns.boxplot(data=df, y='data', x='osrm_distance', palette=cp)
15 plt.title('Data Vs Osrm_distance', fontsize=14, fontfamily='serif', fontweight='bold', backgroundcolor='r', color='w')
16
17 plt.subplot(224)
18 sns.boxplot(data=df, y='route_type', x='actual_distance_to_destination', palette=cp)
19 plt.title('Route_type Vs actual_distance_to_destination', fontsize=14, fontfamily='serif', fontweight='bold', backgroundcolor='k', color=
20 sns.despine()
21 plt.show()

```



Observations:

- Both training and test data have the same range of osrm time recorded
- FTL route type has more actual time compared to Carting. This can also be since FTL is used a lot more than carting in the data available to us
- Both training and test data have the same range of osrm distance recorded
- FTL route type has more actual distance compared to Carting. This can also be since FTL is used a lot more than carting in the data available to us

2. Merging of rows and aggregation of fields

Merging of rows and aggregation of fields

- 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.

```

1 # Grouping by segment
2 # Creating a unique identifier for each segment of a trip
3
4 segment_cols = ['segment_actual_time', 'segment_osrm_distance', 'segment_osrm_time']
5
6 df['segment_key'] = df['trip_uuid'] + ' ' + df['source_center'] + ' ' + df['destination_center']
7
8 for col in segment_cols:
9     df[col + '_sum'] = df.groupby('segment_key')[col].cumsum()
10
11 df[['segment_key', 'segment_actual_time', 'segment_actual_time_sum', 'segment_osrm_distance', 'segment_osrm_distance_sum', 'segment_osr

```

| | segment_key | segment_actual_time | segment_actual_time_sum | segment_osrm_distance | segment_od_start_time |
|--------|---|---------------------|-------------------------|-----------------------|-----------------------|
| 0 | trip-153741093647649320+IND388121AAA+IND388620AAB | 14.0 | 14.0 | 11.965300 | 2023-02-15T10:00:00Z |
| 1 | trip-153741093647649320+IND388121AAA+IND388620AAB | 10.0 | 24.0 | 9.759000 | 2023-02-15T10:00:00Z |
| 2 | trip-153741093647649320+IND388121AAA+IND388620AAB | 16.0 | 40.0 | 10.815200 | 2023-02-15T10:00:00Z |
| 3 | trip-153741093647649320+IND388121AAA+IND388620AAB | 21.0 | 61.0 | 13.022400 | 2023-02-15T10:00:00Z |
| 4 | trip-153741093647649320+IND388121AAA+IND388620AAB | 6.0 | 67.0 | 3.915300 | 2023-02-15T10:00:00Z |
| ... | ... | ... | ... | ... | ... |
| 144862 | trip-153746066843555182+IND131028AAB+IND000000ACB | 12.0 | 92.0 | 8.185800 | 2023-02-15T10:00:00Z |
| 144863 | trip-153746066843555182+IND131028AAB+IND000000ACB | 26.0 | 118.0 | 17.372499 | 2023-02-15T10:00:00Z |
| 144864 | trip-153746066843555182+IND131028AAB+IND000000ACB | 20.0 | 138.0 | 20.705299 | 2023-02-15T10:00:00Z |
| 144865 | trip-153746066843555182+IND131028AAB+IND000000ACB | 17.0 | 155.0 | 18.888500 | 2023-02-15T10:00:00Z |
| 144866 | trip-153746066843555182+IND131028AAB+IND000000ACB | 268.0 | 423.0 | 8.808800 | 2023-02-15T10:00:00Z |

144316 rows × 7 columns

```

1 # Aggregating at segment level & Creating a dictionary for aggregation at segment level
2
3 segment_dict = {
4     'trip_uuid': 'first',
5     'data': 'first',
6     'route_type': 'first',
7     'trip_creation_time': 'first',
8     'source_name': 'first',
9     'destination_name': 'last',
10    'od_start_time': 'first',
11    'od_end_time': 'last',
12    'start_scan_to_end_scan': 'first',
13    'actual_distance_to_destination': 'last',
14    'actual_time': 'last',
15    'osrm_time': 'last',
16    'osrm_distance': 'last',
17    'segment_actual_time': 'sum',
18    'segment_osrm_time' : 'sum',
19    'segment_osrm_distance' : 'sum',
20    'segment_actual_time_sum': 'last',
21    'segment_osrm_time_sum': 'last',
22    'segment_osrm_distance_sum': 'last',
23 }
24
25 # Grouping by segment_key and aggregating
26 segment_agg_data = df.groupby('segment_key').agg(segment_dict).reset_index()
27 segment_agg_data = segment_agg_data.sort_values(by=['segment_key','od_end_time'])
28 segment_agg_data

```

| | segment_key | trip_uuid | data | route_type | trip_creation_time | source |
|-------|---|--------------------|----------|------------|----------------------------|-------------------------------------|
| 0 | trip-153671041653548748+IND209304AAA+IND00000ACB | 153671041653548748 | training | FTL | 2018-09-12 00:00:16.535741 | Kanpur_Central_Uttar Pradesh |
| 1 | trip-153671041653548748+IND462022AAA+IND209304AAA | 153671041653548748 | training | FTL | 2018-09-12 00:00:16.535741 | Bhopal_Transport_Madhya Pradesh |
| 2 | trip-153671042288605164+IND561203AAB+IND562101AAA | 153671042288605164 | training | Carting | 2018-09-12 00:00:22.886430 | Doddablpur_Chikmagalur_Karnataka |
| 3 | trip-153671042288605164+IND572101AAA+IND561203AAB | 153671042288605164 | training | Carting | 2018-09-12 00:00:22.886430 | Tumkur_Veeranahalli_Karnataka |
| 4 | trip-153671043369099517+IND00000ACB+IND160002AAC | 153671043369099517 | training | FTL | 2018-09-12 00:00:33.691250 | Gurgaon_Bilaspur_Haryana |
| ... | ... | ... | ... | ... | ... | ... |
| 26217 | trip-153861115439069069+IND628204AAA+IND627657AAA | 153861115439069069 | test | Carting | 2018-10-03 23:59:14.390954 | Tiruchchendur_Shennongi_Tamil Nadu |
| 26218 | trip-153861115439069069+IND628613AAA+IND627005AAA | 153861115439069069 | test | Carting | 2018-10-03 23:59:14.390954 | Peikulam_SriVnkgudi_Tamil Nadu |
| 26219 | trip-153861115439069069+IND628801AAA+IND628204AAA | 153861115439069069 | test | Carting | 2018-10-03 23:59:14.390954 | Eral_Busstand_Dharmapuri_Tamil Nadu |
| 26220 | trip-153861118270144424+IND583119AAA+IND583101AAA | 153861118270144424 | test | FTL | 2018-10-03 23:59:42.701692 | Sandur_WrdN1C_Karnataka |
| 26221 | trip-153861118270144424+IND583201AAA+IND583119AAA | 153861118270144424 | test | FTL | 2018-10-03 23:59:42.701692 | Hospet_Karnataka |

26222 rows × 20 columns

◆ Understanding:

The rows have been merged based on the unique segment_key, which is a combination of trip_uuid, source_center, and destination_center.

The aggregated dataset reflects the total values for each segment of the trip.

▼ Feature Engineering

```
1 # 1. Calculating time difference between od_start_time and od_end_time
2 segment_agg_data['od_total_time']=(segment_agg_data['od_end_time'] - segment_agg_data['od_start_time'])
3 segment_agg_data['od_time_diff_hour'] = (segment_agg_data['od_total_time']).dt.total_seconds()/3600
4 segment_agg_data
```

| | segment_key | trip_uuid | data | route_type | trip_creation_time | source |
|-------|--|-------------------------|----------|------------|----------------------------|--|
| 0 | 153671041653548748+IND209304AAA+IND00000ACB | trip-153671041653548748 | training | FTL | 2018-09-12 00:00:16.535741 | Kanpur_Central_Uttar Pradesh |
| 1 | 153671041653548748+IND462022AAA+IND209304AAA | trip-153671041653548748 | training | FTL | 2018-09-12 00:00:16.535741 | Bhopal_Transport_Madhya Pradesh |
| 2 | 153671042288605164+IND561203AAB+IND562101AAA | trip-153671042288605164 | training | Carting | 2018-09-12 00:00:22.886430 | Doddaballapur_Chikkaballapur_Karnataka |
| 3 | 153671042288605164+IND572101AAA+IND561203AAB | trip-153671042288605164 | training | Carting | 2018-09-12 00:00:22.886430 | Tumkur_Veeranahalli_Karnataka |
| 4 | 153671043369099517+IND00000ACB+IND160002AAC | trip-153671043369099517 | training | FTL | 2018-09-12 00:00:33.691250 | Gurgaon_Bilaspur_Haryana |
| ... | ... | ... | ... | ... | ... | ... |
| 26217 | 153861115439069069+IND628204AAA+IND627657AAA | trip-153861115439069069 | test | Carting | 2018-10-03 23:59:14.390954 | Tiruchchendur_Shennongi_Tamil Nadu |
| 26218 | 153861115439069069+IND628613AAA+IND627005AAA | trip-153861115439069069 | test | Carting | 2018-10-03 23:59:14.390954 | Peikulam_SriVnkgudi_Tamil Nadu |
| 26219 | 153861115439069069+IND628801AAA+IND628204AAA | trip-153861115439069069 | test | Carting | 2018-10-03 23:59:14.390954 | Eral_Busstand_D |
| 26220 | 153861118270144424+IND583119AAA+IND583101AAA | trip-153861118270144424 | test | FTL | 2018-10-03 23:59:42.701692 | Sandur_WrdN1C_Karnataka |
| 26221 | 153861118270144424+IND583201AAA+IND583119AAA | trip-153861118270144424 | test | FTL | 2018-10-03 23:59:42.701692 | Hospet_Karnataka |

26222 rows × 22 columns

1 segment_agg_data.sample()

| | segment_key | trip_uuid | data | route_type | trip_creation_time | source |
|------|--|-------------------------|----------|------------|----------------------------|--|
| 4965 | 153704666810685062+IND173025AAA+IND160002AAC | trip-153704666810685062 | training | FTL | 2018-09-15 21:24:28.108007 | Paonta Sahib_Gurdwara_Himachal Pradesh |

```
1 # de = segment_agg_data.drop(columns=['source_city', 'source_state', 'source_place', 'destination_place', 'destination_city', 'destination_state'])
2 de = segment_agg_data.copy()
```

1 de.sample()

| | segment_key | trip_uuid | data | route_type | trip_creation_time | source_name |
|------|--|-------------------------|----------|------------|----------------------------|--------------------------------|
| 4962 | 153704666399260818+IND421302AAG+IND401104AAA | trip-153704666399260818 | training | Carting | 2018-09-15 21:24:23.992906 | Bhiwandi_Mankoli_H_Maharashtra |

```
1 # # could have done this --- but some major error ... check it....
2 # sad = segment_agg_data.copy()
3 # sad["source_city"] = sad["source_name"].str.split(" ", n=1, expand=True)[0].str.split("_", n=1, expand=True)[0]
4 # sad["source_state"] = sad["source_name"].str.split(" ", n=1, expand=True)[1].str.replace("(", "").str.replace(")", "")
5
6 # sad["destination_city"] = sad["destination_name"].str.split(" ", n=1, expand=True)[0].str.split("_", n=1, expand=True)[0]
7 # sad["destination_state"] = sad["destination_name"].str.split(" ", n=1, expand=True)[1].str.replace("(", "").str.replace(")", "")
8
9 # sad["source_place"] = sad["source_name"].str.split("_", n=2, expand=True)[1]
10 # sad["destination_place"] = sad["destination_name"].str.split("_", n=2, expand=True)[1]
```

```
1 # using regex pattern to separate the city, place, state
2 def extract_info(name):
3     pattern = r'^^(?P<city>[^\\s_]+)_?(?P<place>[^\\(\\)]*)\\s?\\((?P<state>[A-Za-z\\s&]+)\\)$'
4     match = re.match(pattern, name)
5     if match:
6         city = match.group('city').strip()
7         place = match.group('place').strip() if match.group('place') else city
8         state = match.group('state').strip()
9         return city, place, state
10    else:
11        return None, None, None
```

1 de[['source_city', 'source_place', 'source_state']] = de['source_name'].apply(lambda x: pd.Series(extract_info(x)))

```
1 de[['destination_city', 'destination_place', 'destination_state']] = de['destination_name'].apply(lambda x: pd.Series(extract_info(x))
```

```
1 de
```

| | segment_key | trip_uuid | data | route_type | trip_creation_time | source |
|-------|---|--------------------|----------|------------|----------------------------|--|
| 0 | trip-153671041653548748+IND209304AAA+IND00000ACB | 153671041653548748 | training | FTL | 2018-09-12 00:00:16.535741 | Kanpur_Central (Uttar Pradesh) |
| 1 | trip-153671041653548748+IND462022AAA+IND209304AAA | 153671041653548748 | training | FTL | 2018-09-12 00:00:16.535741 | Bhopal_Transport (Madhya Pradesh) |
| 2 | trip-153671042288605164+IND561203AAB+IND562101AAA | 153671042288605164 | training | Carting | 2018-09-12 00:00:22.886430 | Doddaballapur_Chikkaballapur (Karnataka) |
| 3 | trip-153671042288605164+IND572101AAA+IND561203AAB | 153671042288605164 | training | Carting | 2018-09-12 00:00:22.886430 | Tumkur_Veerapuram (Karnataka) |
| 4 | trip-153671043369099517+IND00000ACB+IND160002AAC | 153671043369099517 | training | FTL | 2018-09-12 00:00:33.691250 | Gurgaon_Bilaspur (Haryana) |
| ... | ... | ... | ... | ... | ... | ... |
| 26217 | trip-153861115439069069+IND628204AAA+IND627657AAA | 153861115439069069 | test | Carting | 2018-10-03 23:59:14.390954 | Tiruchchendur_Shennongi (Tamil Nadu) |
| 26218 | trip-153861115439069069+IND628613AAA+IND627005AAA | 153861115439069069 | test | Carting | 2018-10-03 23:59:14.390954 | Peikulam_SriVnkti (Tamil Nadu) |
| 26219 | trip-153861115439069069+IND628801AAA+IND628204AAA | 153861115439069069 | test | Carting | 2018-10-03 23:59:14.390954 | Eral_Busstand_D |
| 26220 | trip-153861118270144424+IND583119AAA+IND583101AAA | 153861118270144424 | test | FTL | 2018-10-03 23:59:42.701692 | Sandur_WrdN1D (Karnataka) |
| 26221 | trip-153861118270144424+IND583201AAA+IND583119AAA | 153861118270144424 | test | FTL | 2018-10-03 23:59:42.701692 | Hospet (Karnataka) |

26222 rows × 28 columns

◀ ▶

```
1 de[(de['source_place']=='') | (de['destination_place']=='')]
```

| | segment_key | trip_uuid | data | route_type | trip_creation_time | source_name |
|-------|---|--------------------|----------|------------|----------------------------|------------------------------------|
| 7 | trip-153671052974046625+IND583101AAA+IND583201AAA | 153671052974046625 | training | FTL | 2018-09-12 00:02:09.740725 | Bellary_District (Karnataka) |
| 9 | trip-153671052974046625+IND583201AAA+IND583119AAA | 153671052974046625 | training | FTL | 2018-09-12 00:02:09.740725 | Hospet (Karnataka) |
| 19 | trip-153671110078355292+IND121004AAB+IND121001AAA | 153671110078355292 | training | Carting | 2018-09-12 00:11:40.783923 | FBD_Balabhgarh_District (Haryana) |
| 33 | trip-153671173668736946+IND110043AAA+IND110078AAA | 153671173668736946 | training | Carting | 2018-09-12 00:22:16.687619 | Delhi_Nangli (Delhi) |
| 80 | trip-153671320807895983+IND121004AAB+IND121102AAA | 153671320807895983 | training | Carting | 2018-09-12 00:46:48.079257 | FBD_Balabhgarh_District (Haryana) |
| ... | ... | ... | ... | ... | ... | ... |
| 26118 | trip-153860849934816308+IND110078AAA+IND110043AAA | 153860849934816308 | test | Carting | 2018-10-03 23:14:59.348414 | Janakpuri (Delhi) |
| 26153 | trip-153860958923357924+IND842003AAB+IND482002AAA | 153860958923357924 | test | Carting | 2018-10-03 23:33:09.233829 | Jabalpur_Adhartal (Madhya Pradesh) |
| 26180 | trip-153861007249500192+IND842001AAA+IND846004AAA | 153861007249500192 | test | FTL | 2018-10-03 23:41:12.495257 | Muzaffarpur_Bihar (Bihar) |
| 26181 | trip-153861007249500192+IND846004AAA+IND847103AAA | 153861007249500192 | test | FTL | 2018-10-03 23:41:12.495257 | Darbhanga (Bihar) |
| 26221 | trip-153861118270144424+IND583201AAA+IND583119AAA | 153861118270144424 | test | FTL | 2018-10-03 23:59:42.701692 | Hospet (Karnataka) |

782 rows × 28 columns

◀ ▶

```
1 de.loc[de['source_place']=='','source_place']=de['source_city']
2 de.loc[de['destination_place']=='','destination_place']=de['destination_city']
```

```
1 de[de.source_place.isna()]
```

| | segment_key | trip_uuid | data | route_type | trip_creation_time | source_name | destination_name | od_start_time | od_end_time | start_sca |
|--|-------------|-----------|------|------------|--------------------|-------------|------------------|---------------|-------------|-----------|
| | | | | | | | | | | |

```
1 de.isna().sum()

[segment_key, trip_uuid, data, route_type, trip_creation_time, source_name, destination_name, od_start_time, od_end_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, segment_actual_time_sum, segment_osrm_time_sum, segment_osrm_distance_sum, od_total_time, od_time_diff_hour, source_city, source_place, source_state, destination_city, destination_place, destination_state]
dtype: int64

1 #de = de.drop(columns=['source_city','source_state','source_place','destination_city', 'destination_place', 'destination_state'])

1 de.loc[de.source_city=='Bangalore','source_city']='Bengaluru'
2 de.loc[de.destination_city=='Bangalore','destination_city']='Bengaluru'

1 np.set_printoptions(threshold=np.inf)

1 de['source_city'].unique()
```

```

'Kannanakul', 'Kothanapattioor', 'Nirjuvai', 'Siunmukri', 'Kuttiu',
'Manteswar', 'Ranaghat', 'Gudalur', 'Kasganj', 'Helencha',
'Dabhoi', 'Balangir', 'Bargarh', 'Sirohi', 'Gondal', 'Morena',
'Nagpur', 'Baripada', 'Palitana', 'Samsi', 'Sumerpur', 'Vadipatti',
'Giddarbarha', 'Fazilka', 'Phusro', 'Tonk', 'SawaiMadhopur',
'Munger', 'Pauri', 'Sankerko', 'Nohar', 'Krishnanagar',
'Kuthuparamba', 'KharagpurBR', 'Jamui', 'DhrmpuriTS', 'Anuppur',
'Baruipur', 'Tirtol', 'Daurala', 'Jharsuguda', 'Central', 'Abohar',
'RampuraPhul', 'Hoshangabad', 'Cuttack', 'Kendrpara', 'Rajgarh',
'Durg', 'Balurghat', 'Chamorshi', 'Barasat', 'Gujilam', 'Basta',
'Sahjhnpur', 'Tadipatri', 'Cuddapah', 'Badvel', 'Proddatur',
'Angul', 'Phulbani', 'Ambegaon', 'Chanchal', 'Malout',
'Uthangarai', 'Badlapur', 'Balugaon', 'Mahasamund', 'Badkulla',
'Kanadwanii' ], dtype=object)

```

```
1 de['source_place'].unique()
```

```

[{"HotelPrk_D", "SnthiNGR_D", "BypRDDPP_D", "BhmrDPP_D",
 "Vijayawada", "Kidwai_D", "Mharajpr_D", "GrudwrD_D", "Adrshngr_D",
 "Krvnkuzy_D", "PhdofDPP_D", "CtyLgDPP_D", "Pilani", "StnRoad_DC",
 "Solaiprm_D", "RailwyRd_D", "Diakkawn_D", "PiliKoti_D",
 "VislkNgr_DC", "SaiTemp1_D", "NcsRd_DC", "Arangadi_D",
 "LohiaDPP_D", "PakriDPP_D", "BisnoiDPP_D", "PrmnRDP_D",
 "BsnoiHPL_D", "Gopa3PL_D", "ShrprDPP_D", "MnRDP_D", "BllvMarg_D",
 "Farmnala_D", "Katghara_D", "KhandDPP_D", "JngidDPP_D",
 "Tolichwk_L", "Pringla_D", "SurbhiTh_D", "Greens_D", "CivilStn_D",
 "Shivprsd_D", "KirtinGnr_D", "ModelTwn_P", "Fairybnk_D",
 "Cnsrvila_D", "BhunaDPP_D", "Ukkadam_D", "BChwkDPP_D", "Jharia_DC",
 "Chrwpaty_D", "CollgerD_D", "Pthrgoan_D", "HelipadRD_D",
 "DlptrDPP_D", "Ranakant_D", "Vandalur_Dc", "Airport",
 "Rahatani_DPC", "Bhrnikvu_D", "HunthrVg_I", "Awmpivng_D",
 "PmthuKlm_D", "PanditRd_D", "Chtprpuza_D", "SantaNGR_D",
 "JalnaRd_D", "Bazzard_D", "UmarDPP_D", "CmtNgRod_D", "KasyaDPP_D",
 "StteHW28_D", "VijywdRD_D", "Dehrird_D", "Wardno7_D", "NatunDPP_D",
 "Enayetpr_D", "East_I_20", "KtnRdDPP_D", "Mapusa", "Kothriya_DC",
 "Ymunpurm_D", "TonkRoad_D", "East_D_8", "SmbjiCwk_D", "Jabalpur",
 "Indrapri_DC", "YashDPP_D", "AchneraRD_D", "DibngVly_D",
 "BnglorRd_D", "PcrndDPP_D", "Murtngr_D", "Central_D_15",
 "JmpuCntn_D", "Rajpura_D", "Mulapali_D", "LdnunDPP_D", "PngnrRd_D",
 "Tejpal_M", "RjghatRd_D", "Pnchlght_D", "KDRoad_D", "farukngr_D",
 "SingCLNY_D", "Jagrata_D", "ThaneDPP_D", "MdhsnPpp_D", "Margao_Dc",
 "Nijgan_D", "Rawatpur_D", "Gurukrpa_D", "JivanDPP_D",
 "AshokNagar_DC", "BhmrDPP_D", "KarjuDPP_D", "Pothredy_D", "_NAD",
 "ShbdnDPP_D", "ClgRDDPP_D", "Brpc", "ShubsNGR_D", "Beltnngdi_D",
 "Ldthlabh_D", "KamaDPP_D", "InDestat_L", "HnmntNgr_D",
 "ShivmDPP_D", "BstndDPP_D", "Cherukole_D", "EmsPnmbi_D",
 "KaimgnjRD_D", "Bhainpura_D", "FatprDPP_D", "RajRdDPP_D",
 "Mahad_D", "MsjidDPP_D", "PaikjNGR_D", "BalibDPP_D", "KcharaRD_D",
 "ChtwrDPP_D", "MsmcyDPP_D", "ElngoNgr_C", "Rcocmplx_D",
 "PunjbiPd_D", "Shanthi_D", "Ponda_Dc", "SirsaDPP_D", "Mahim",
 "Paldi_D", "SubrtDPP_D", "Santalpr_D", "TiloDPP_D", "Idgah_L",
 "KoralDPP_D", "Nerul_D", "Uppal_Dc", "MhliaDPP_D", "Bomsndra_L",
 "SukntDPP_D", "Arsprmbu_D", "MnmlaRd_D", "UttarDPP_D",
 "ChatidPP_D", "MdnpDPP_D", "Kanakpur_D", "East", "MdothdRD_D",
 "CroadDPP_D", "Sarjapur_D", "SbhRDDP_D", "Central_D_4",
 "Thiruvlr_DC", "Sangetha_D", "NadthiCx_D", "Naraynpr_D", "Hathras",
 "AkkolRD_D", "AnnaNGR_D", "LamtidPP_D", "DiyoDPP_D", "MunplDPP_D",
 "RoopNgr_D", "RwStnDPP_D", "BhrolDPP_D", "Umargaon_DC",
 "Rammagar_D", "Majoor_D", "BageDPP_D", "MnbzrDPP_D", "AkhraBzr_D",
 "MemariRD_D", "ArickDPP_D", "kalmpuza_D", "BnkrGate_D",
 "ColnyDPP_D", "Muktars_D", "Patia", "Rjndpara_D", "BreezeDPP_D",
 "GurpdDPP_D", "Kothanur_L", "Ricco_D", "Palam", "Prjapati_D",
 "Gondkhry_H", "South_D_20", "Central_H_4", "KalikDPP_D",
 "Bareilly", "STRdDPP_D", "RatuaDPP_D", "BazarDPP_D", "lalaNGR_D",
 "GndhiChk_D", "Mirapati_L", "RhmgjDPP_D", "Central_L_8", "Kaura_D",
 "Khndyusn_D", "AmbkaDPP_D", "NeImngla_L", "IOTCEnl_L",
 "AnadiDPP_D", "JJCpxDPP_D", "IdstrlAr_D", "Mehsouri_D",
 "HanumDPP_D", "GovndNgr_DC", "JthridPP_DC", "Karelibaug_DC",
 "Bnsibtla_D", "MjlprDPP_D", "Sardhnrd_D", "MbRoad_D",
 "Central_D_10", "Old", "SliprDPP_DC", "Kapleswr_D", "Kdvantra_D",
 "SadulDPP_D", "Pulgano_DC", "AkhirDPP_D", "KrsnNgr_D", "Parai_D",
 "Selaiyur_DC", "Chinchwad DC", "NavdaCln_D", "CBRoad_D", "Sidrd_D",
 "ThanaDPP_D", "Kruspharma_D", "Mahuva_DC", "Manchar_D",
 "BargaDPP_D", "NapitDPP_D", "RgstrOFC_D", "Bhaluahi_D",
 "Thikiri_D", "RajpurRD_D"], dtype=object)

```

```
1 de['source_state'].unique()
```

```

[{"array(['Uttar Pradesh', 'Madhya Pradesh', 'Karnataka', 'Haryana',
 'Maharashtra', 'Tamil Nadu', 'Gujarat', 'Delhi', 'Telangana',
 'Andhra Pradesh', 'Rajasthan', 'Assam', 'West Bengal', 'Punjab',
 'Chandigarh', 'Goa', 'Uttarakhand', 'Jharkhand', 'Pondicherry',
 'Orissa', 'Himachal Pradesh', 'Kerala', 'Arunachal Pradesh',
 'Bihar', 'Meghalaya', 'Chhattisgarh', 'Jammu & Kashmir',
 'Dadra and Nagar Haveli', 'Mizoram', 'Tripura', 'Nagaland'],
 dtype=object)}

```

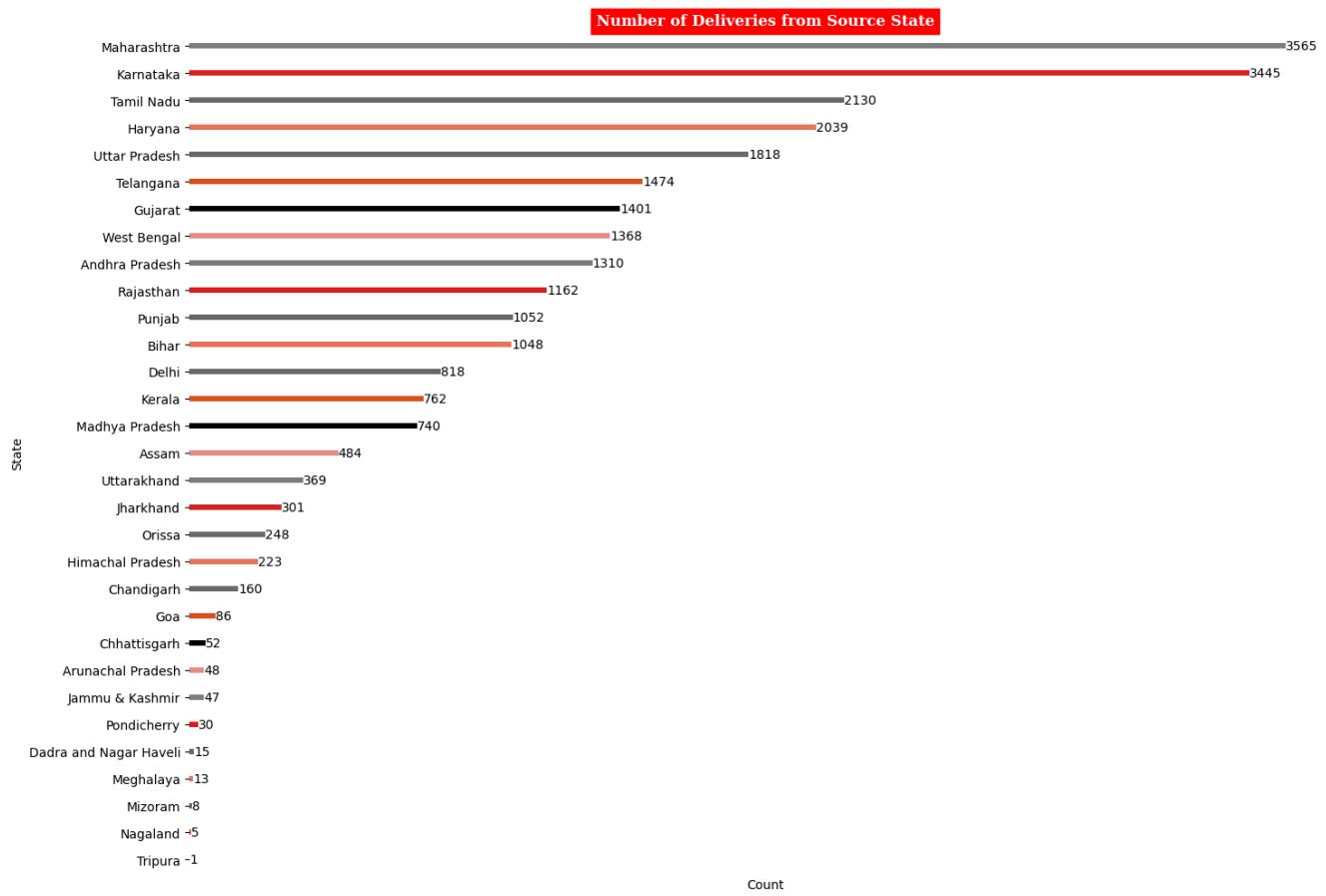
```
1 de['source_state'].value_counts().to_frame().style.background_gradient(cmap='Reds')
```

| | count |
|-------------------------------|-------|
| source_state | |
| Maharashtra | 3565 |
| Karnataka | 3445 |
| Tamil Nadu | 2130 |
| Haryana | 2039 |
| Uttar Pradesh | 1818 |
| Telangana | 1474 |
| Gujarat | 1401 |
| West Bengal | 1368 |
| Andhra Pradesh | 1310 |
| Rajasthan | 1162 |
| Punjab | 1052 |
| Bihar | 1048 |
| Delhi | 818 |
| Kerala | 762 |
| Madhya Pradesh | 740 |
| Assam | 484 |
| Uttarakhand | 369 |
| Jharkhand | 301 |
| Orissa | 248 |
| Himachal Pradesh | 223 |
| Chandigarh | 160 |
| Goa | 86 |
| Chhattisgarh | 52 |
| Arunachal Pradesh | 48 |
| Jammu & Kashmir | 47 |
| Pondicherry | 30 |
| Dadra and Nagar Haveli | 15 |
| Meghalaya | 13 |
| Mizoram | 8 |
| Nagaland | 5 |
| Tripura | 1 |

```

1 state_counts = de['source_state'].value_counts().to_frame().reset_index()
2 state_counts.columns = ['State', 'Count']
3
4 plt.figure(figsize=(15,10))
5 a = sns.barplot(y='State', x='Count', data=state_counts, palette=cp, width=0.2)
6 a.bar_label(a.containers[0], label_type='edge')
7 plt.xticks([])
8 plt.ylabel('State')
9 plt.xlabel('Count')
10 plt.title('Number of Deliveries from Source State', fontsize=12, fontfamily='serif', fontweight='bold', backgroundcolor='r', color='w')
11 plt.tight_layout()
12 sns.despine(bottom=True, left=True)
13 plt.show()

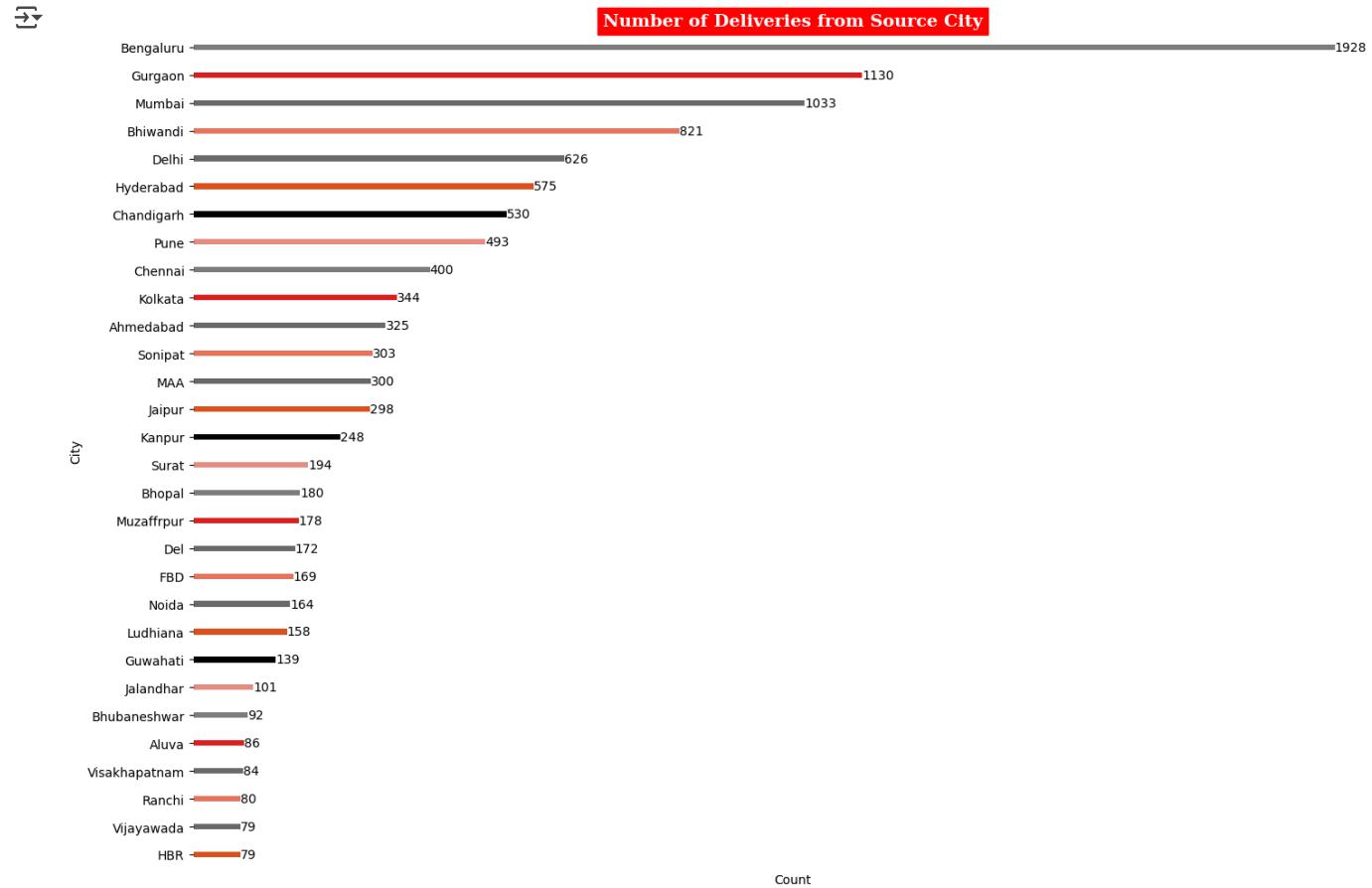
```



```

1 city_counts = de['source_city'].value_counts().to_frame().reset_index()[:30]
2 city_counts.columns = ['City', 'Count']
3
4 plt.figure(figsize=(15,10))
5 a = sns.barplot(y='City', x='Count', data=city_counts, palette=cp, width=0.2)
6 a.bar_label(a.containers[0], label_type='edge')
7 plt.xticks([])
8 plt.ylabel('City')
9 plt.xlabel('Count')
10 plt.title('Number of Deliveries from Source City', fontsize=14, fontfamily='serif', fontweight='bold', backgroundcolor='r', color='w')
11 plt.tight_layout()
12 sns.despine(bottom=True, left=True)
13 plt.show()

```



⌄ ⚡ Insights:

Source State contributors

- **Maharashtra, Karnataka, Tamil Nadu, Haryana, and Uttar Pradesh** are the top contributors where maximum bookings are recorded in this month indicating significant engagement.

Source City contributors

- Cities like **Bengaluru, Gurgaon, Mumbai, Bhiwandi, Delhi, Hyderabad** where the major no.of booking are recorded.

```
1 de.describe().T
```

| | count | mean | min | 25% | 50% | 75% |
|--------------------------------|---------|----------------------------------|-------------------------------|----------------------------------|----------------------------------|--|
| trip_creation_time | 26222 | 2018-09-22 13:58:56.740969728 | 2018-09-12 00:00:16.535741 | 2018-09-17 03:57:50.900417024 | 2018-09-22 03:33:30.255023104 | 2018-09-27 19:55:17.273207040 23:5 |
| od_start_time | 26222 | 2018-09-22 17:49:54.012840448 | 2018-09-12 00:00:16.535741 | 2018-09-17 07:43:36.525784320 | 2018-09-22 07:17:15.379571712 | 2018-09-27 23:22:20.105725952 04:2 |
| od_end_time | 26222 | 2018-09-22 22:49:00.449498112 | 2018-09-12 00:50:10.814399 | 2018-09-17 15:10:35.615369472 | 2018-09-22 14:46:05.410478848 | 2018-09-28 03:19:50.211971840 03:0 |
| start_scan_to_end_scan | 26222.0 | 298.553375 | 20.0 | 90.0 | 152.0 | 307.0 |
| actual_distance_to_destination | 26222.0 | 92.533051 | 9.001351 | 21.65415 | 35.044329 | 65.557392 1 |
| actual_time | 26222.0 | 200.92659 | 9.0 | 51.0 | 84.0 | 167.0 |
| osrm_time | 26222.0 | 90.785332 | 6.0 | 25.0 | 39.0 | 72.0 |
| osrm_distance | 26222.0 | 114.975334 | 9.0729 | 27.71915 | 43.54355 | 85.443949 2 |
| segment_actual_time | 26222.0 | 199.095642 | 9.0 | 50.0 | 83.0 | 166.0 |
| segment_osrm_time | 26222.0 | 101.793343 | 6.0 | 25.0 | 42.0 | 79.0 |
| segment_osrm_distance | 26222.0 | 125.587128 | 9.0729 | 28.429099 | 45.797649 | 91.023573 2 |
| segment_actual_time_sum | 26222.0 | 199.095642 | 9.0 | 50.0 | 83.0 | 166.0 |
| segment_osrm_time_sum | 26222.0 | 101.793343 | 6.0 | 25.0 | 42.0 | 79.0 |
| segment_osrm_distance_sum | 26222.0 | 125.587128 | 9.0729 | 28.429099 | 45.797649 | 91.023573 2 |
| od_total_time | 26222 | 0 days 04:59:06.436657725 | 0 days 00:20:42.168787 | 0 days 01:30:58.665517750 | 0 days 02:32:20.203949500 | 0 days 05:07:17.158769 11:3 |
| od_time_diff_hour | 26222.0 | 4.985121 | 0.345047 | 1.516296 | 2.538946 | 5.121433 |

```
1 de.describe(include='object').T
```

| | count | unique | top | freq |
|-------------------|-------|---|------|------|
| segment_key | 26222 | 26222 trip-153671041653548748+IND209304AAA+IND000000ACB | 1 | |
| trip_uuid | 26222 | 14787 trip-153717306559016761 | 8 | |
| source_name | 26222 | 1496 Gurgaon_Bilaspur_HB (Haryana) | 1052 | |
| destination_name | 26222 | 1466 Gurgaon_Bilaspur_HB (Haryana) | 928 | |
| source_city | 26222 | 1239 Bengaluru | 1928 | |
| source_place | 26222 | 1246 Bilaspur_HB | 1052 | |
| source_state | 26222 | 31 Maharashtra | 3565 | |
| destination_city | 26222 | 1236 Bengaluru | 1863 | |
| destination_place | 26222 | 1217 Bilaspur_HB | 928 | |
| destination_state | 26222 | 32 Karnataka | 3497 | |

```
1 de['destination_city'].unique()
```

Bikramgarh , Bellampalli , Raigarh , Kaliyaganj , Parwanpur ,
'Sihora', 'Shivpuri', 'Gondal', 'Khurja', 'Tonk', 'Parbhani',
'Suratgarh', 'Mandsaur', 'Sitamau', 'Agra', 'Amalner', 'Pachora',
'Erandoi', 'Pataudi', 'LakhimpurN', 'Gopiganj', 'Baripada',
'Mughalsarai', 'Palamaner', 'Umerkote', 'Shamshabad', 'Dhuri',
'Farrukhbad', 'Shahabad', 'Shahjhpur', 'Aliganj', 'Jalalabad',
'Palakollu', 'Samsi', 'Wankaner', 'BariSadri', 'Pilibanga',
'Vishakhapatnam', 'Lodhan', 'Hoshangabad', 'Bokakhat', 'Moranhat',
'Anjar', 'Barauni', 'Gahmar', 'Jowai', 'Varanasi', 'Chiraiyakot',
'ChrkhiDdri', 'Kozhenchery', 'Chhibramau', 'Karukachal',
'Parakkadavu', 'Nimbahera', 'Neemuch', 'MirzapurWB', 'Berhampore',
'Manbazar', 'Khatra', 'Kothanalloor', 'Badlapur', 'Saraimeer',
'Shahganj', 'Kerala', 'Buhana', 'Tallada', 'Itahan', 'Manikchak',
'Sonari', 'Balotra', 'Simlapal', 'Ghatkesar', 'Helencha',
'Jaleswar', 'Soro', 'Manglamaro', 'Keshiary', 'Contai', 'Howrah',
'Udala', 'Chikmagalur', 'Kasganj', 'Umaria', 'MarwarJn', 'Jalore',
'Thirthahalli', 'Khanakul', 'Chikodi', 'Khamgaon', 'Sundargarh',
'Rajgangpur', 'Bargarh', 'Bhinmal', 'Sirohi', 'Dabhoi',
'RampuraPhul', 'Bhilad', 'Arambag', 'Vaikom', 'Manteswar',
'Nabadwip', 'Gangarampr', 'Bongaon', 'Morena', 'Peddapalli',
'Roing', 'Fazilka', 'Cjb', 'Sonepur', 'Sambalpur', 'Balangir',
'Sumerpur', 'Jasdan', 'Khalilabad', 'Nagpur', 'Aunrihar',
'Mungaoli', 'Chanchal', 'Falna', 'Gujiliam', 'Malout', 'Gomoh',
'Kamarpukur', 'AmaDubi', 'Azamgarh', 'Phusro', 'Tirurangadi',
'SawaiMadhopur', 'Sultanganj', 'Giddarbarha', 'Bhadra', 'Hanskhali',
'Munger', 'KharagpurBR', 'Baruipur', 'Khatauli', 'Ambegaon',
'Abohar', 'Nirsa', 'Itarsi', 'Cuttack', 'Taranagar', 'Kanker',
'Luxettipet', 'Basta', 'Bhadohi', 'Tilhar', 'Tadipatri',
'Proddatur', 'Badvel', 'Cuddapah', 'Angul', 'Phulbani', 'Khetri',
'Daman', 'Uthangarai', 'Dung', 'Thachnttukra', 'Krishnanagar',
'Ranaghat', 'Kapadvanj', 'Lunawada'], dtype=object)

```
1 de['destination_place'].unique()
```



```
[Sanpada_CP', 'Bnilai_DC', 'Nattukal_D', 'Anadippp_D',
'ArickDPP_D', 'VrdhriRD_D'], dtype=object)
```

```
1 de['destination_state'].unique()
```

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

```
1 state_counts = de['destination_state'].value_counts().to_frame().reset_index()
```

```
2 state_counts.columns = ['State', 'Count']
```

```
3
```

```
4 plt.figure(figsize=(15,10))
```

```
5 a = sns.barplot(y='State', x='Count', data=state_counts, palette=cp, width=0.2)
```

```
6 a.bar_label(a.containers[0], label_type='edge')
```

```
7 plt.xticks([])
```

```
8 plt.ylabel('State')
```

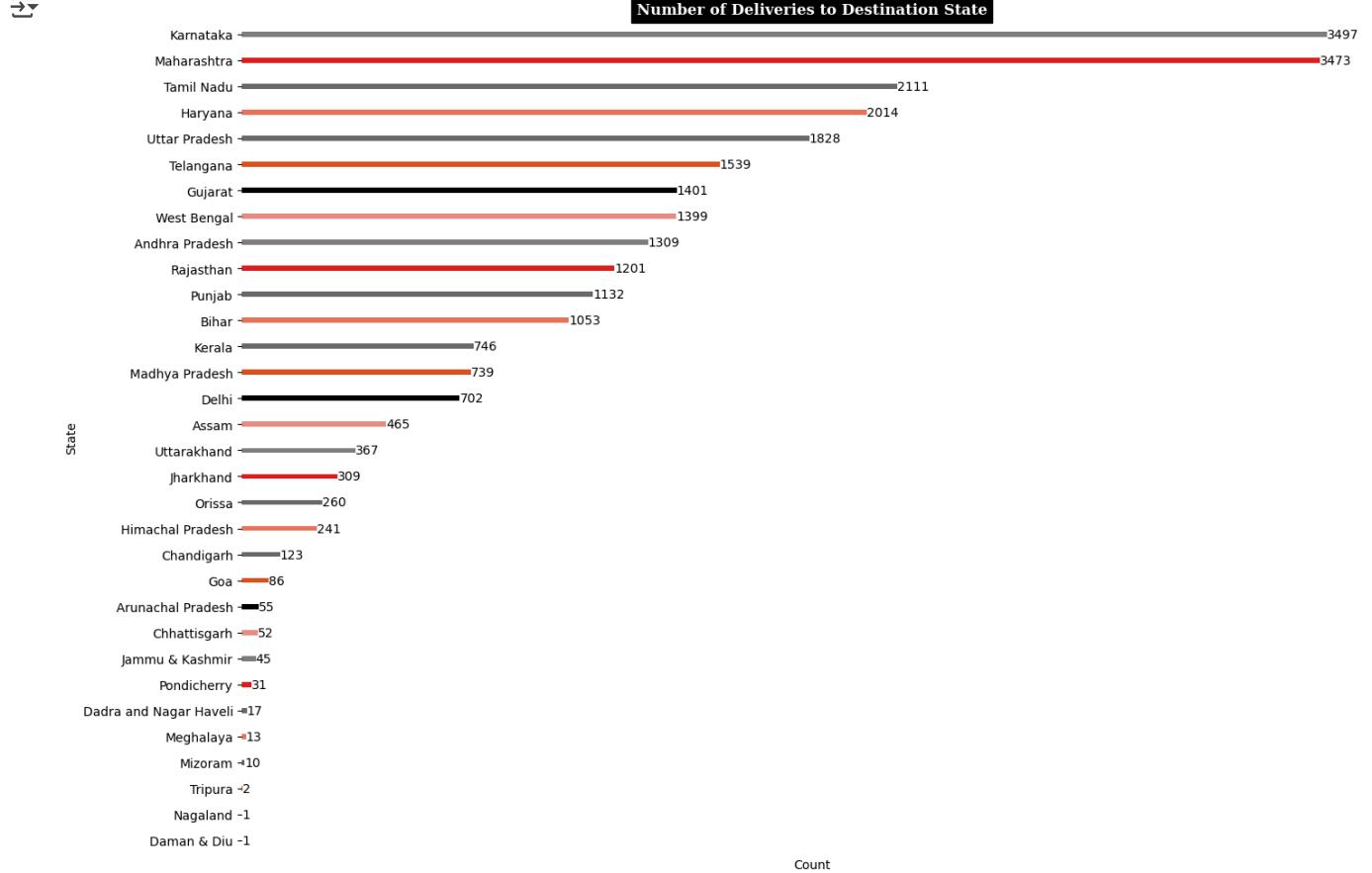
```
9 plt.xlabel('Count')
```

```
10 plt.title('Number of Deliveries to Destination State', fontsize=12, fontfamily='serif', fontweight='bold', backgroundcolor='k', color='w')
```

```
11 plt.tight_layout()
```

```
12 sns.despine(bottom=True, left=True)
```

```
13 plt.show()
```



```
1 city_counts = de['destination_city'].value_counts().to_frame().reset_index()[:30]
```

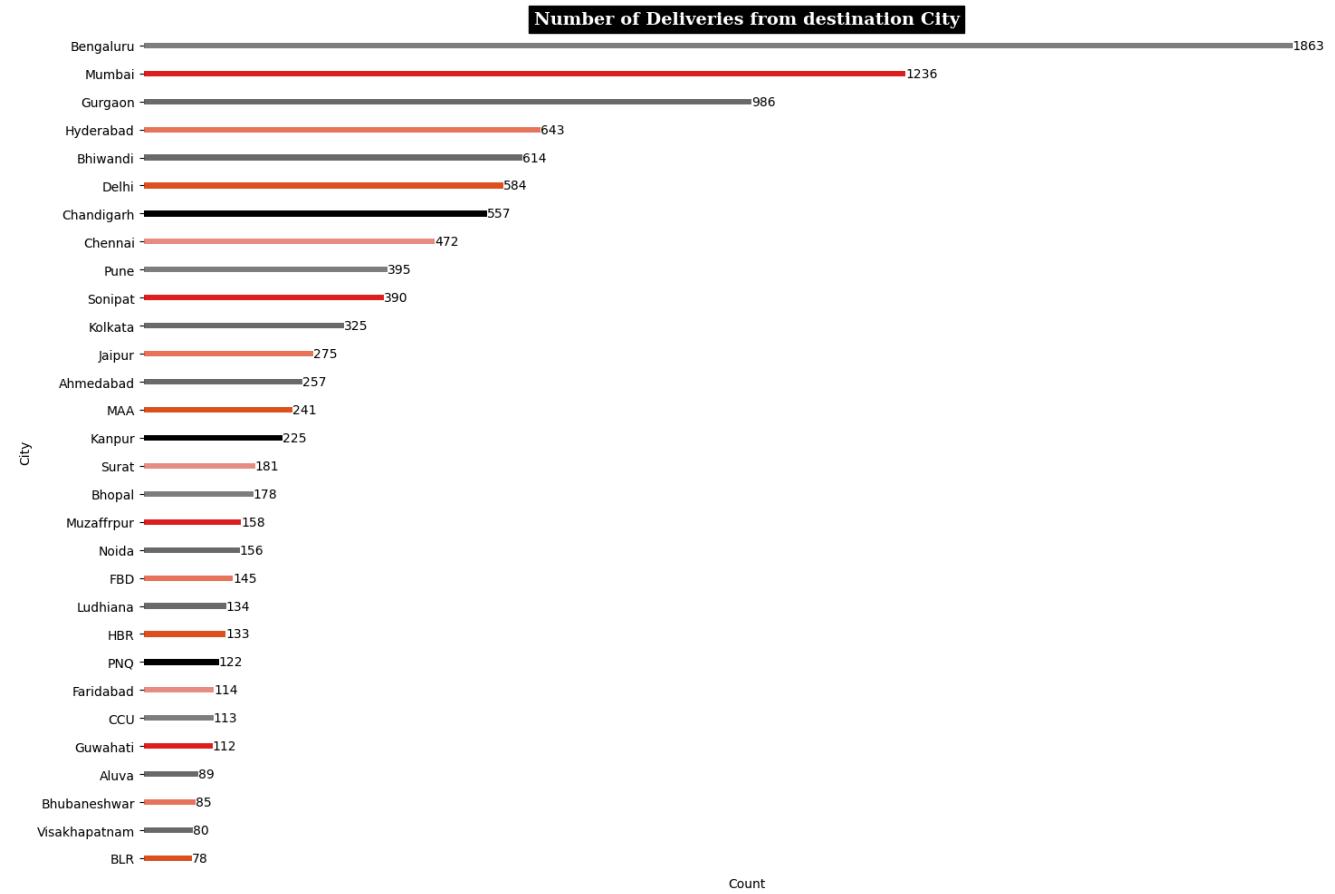
```
2 city_counts.columns = ['City', 'Count']
```

```
3
```

```

4 plt.figure(figsize=(15,10))
5 a = sns.barplot(y='City', x='Count', data=city_counts, palette=cp, width=0.2)
6 a.bar_label(a.containers[0], label_type='edge')
7 plt.xticks([])
8 plt.ylabel('City')
9 plt.xlabel('Count')
10 plt.title('Number of Deliveries from destination City', fontsize=14, fontfamily='serif', fontweight='bold', backgroundcolor='k', color='w')
11 plt.tight_layout()
12 sns.despine(bottom=True, left=True)
13 plt.show()

```



▼ 🌐 Insights:

Destination State

- States like **Karnataka, Maharashtra, Tamil Nadu, Haryana, and Uttar Pradesh** where maximum packages are received in this month indicating significant engagement.

Destination City

- Cities like **Bengaluru, Mumbai, Gurgaon, Bhiwandi, Hyderabad, Delhi** where the major no.of booking are received.

```
1 np.set_printoptions(threshold=np.inf)
```

```

1 de['corridor'] = de['source_name'] + ' <--> ' + de['destination_name']
2 de['corridor'].value_counts()

```

```

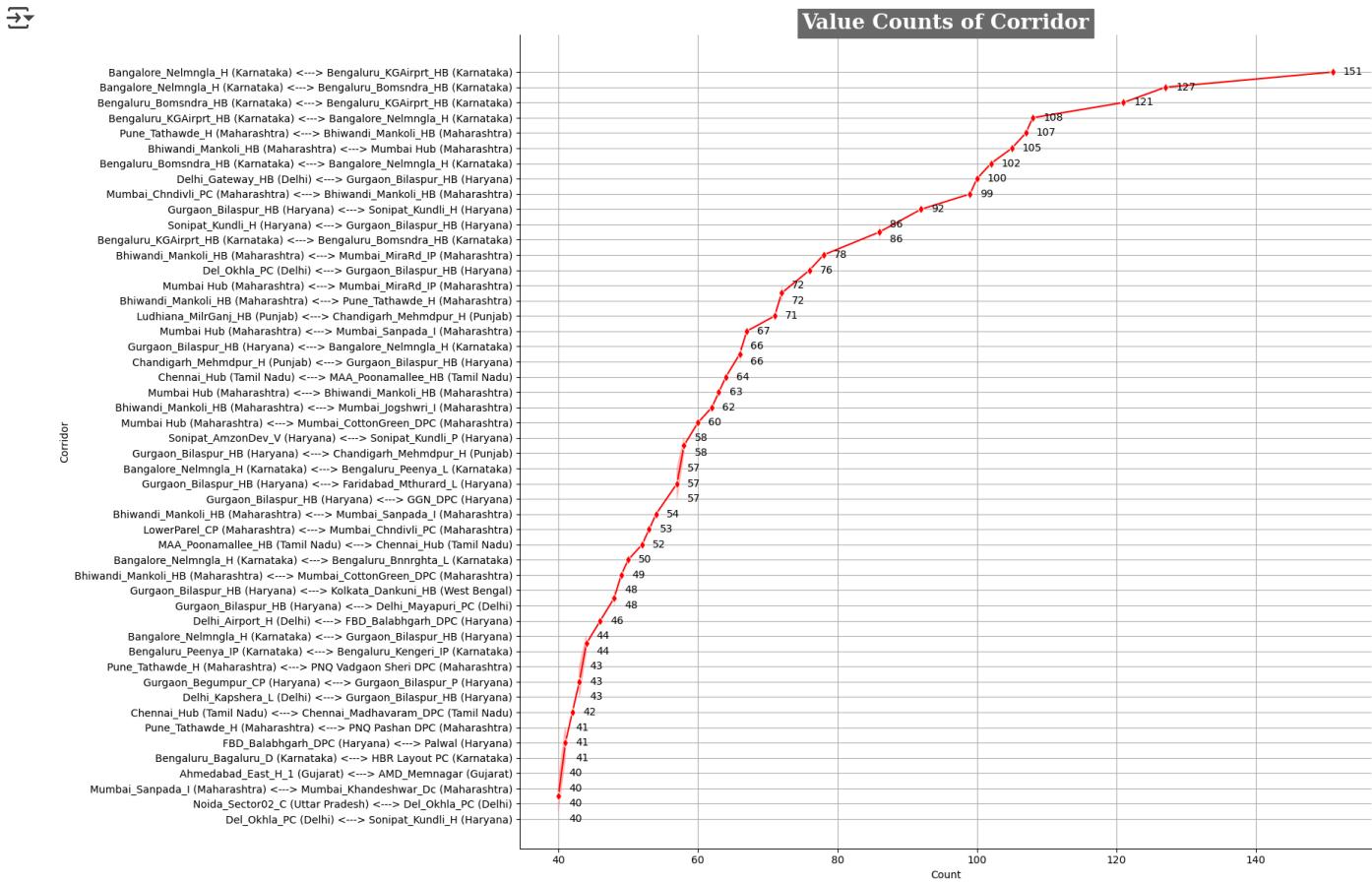
corridor
Bangalore_Nelmgla_H (Karnataka) <--> Bengaluru_KGAirprt_HB (Karnataka) 151
Bangalore_Nelmgla_H (Karnataka) <--> Bengaluru_Bomsndra_HB (Karnataka) 127
Bengaluru_Bomsndra_HB (Karnataka) <--> Bengaluru_KGAirprt_HB (Karnataka) 121
Bengaluru_KGAirprt_HB (Karnataka) <--> Bangalore_Nelmgla_H (Karnataka) 108
Pune_Tathawde_H (Maharashtra) <--> Bhiwandi_Mankoli_HB (Maharashtra) 107
...
Ongole_SubhVRTL_I (Andhra Pradesh) <--> Kandukur_LICOffce_D (Andhra Pradesh) 1
Madnapalle_PngnrRd_D (Andhra Pradesh) <--> Palamaner_Lakshmi_D (Andhra Pradesh) 1
Dharmavram_SaiNgr_D (Andhra Pradesh) <--> Kadiri_GVManu_D (Andhra Pradesh) 1
Baharampur_Chuanpur_I (West Bengal) <--> Chapra_NagarDPP_D (West Bengal) 1
Jaipur_NgrNigam_DC (Rajasthan) <--> Jaipur_Central_D_1 (Rajasthan) 1
Name: count, Length: 2741, dtype: int64

```

```

1 corridor_counts = de['corridor'].value_counts()[:50]
2
3 plt.figure(figsize=(18,12))
4 #corridor_counts.plot(kind='line', marker='d', color='r')
5 sns.lineplot(y=corridor_counts.index, x=corridor_counts.values, marker='d', color='r')
6 plt.title('Value Counts of Corridor', fontsize=20, fontfamily='serif', fontweight='bold', backgroundcolor='dimgray', color='w')
7 plt.ylabel('Corridor')
8 plt.xlabel('Count')
9 plt.tight_layout()
10 sns.despine()
11 plt.grid(True)
12
13 for i, count in enumerate(corridor_counts.values):
14     plt.text(count+1.5, corridor_counts.index[i], str(count), ha='left', va='center')
15
16 plt.show()

```



❖ Insights:

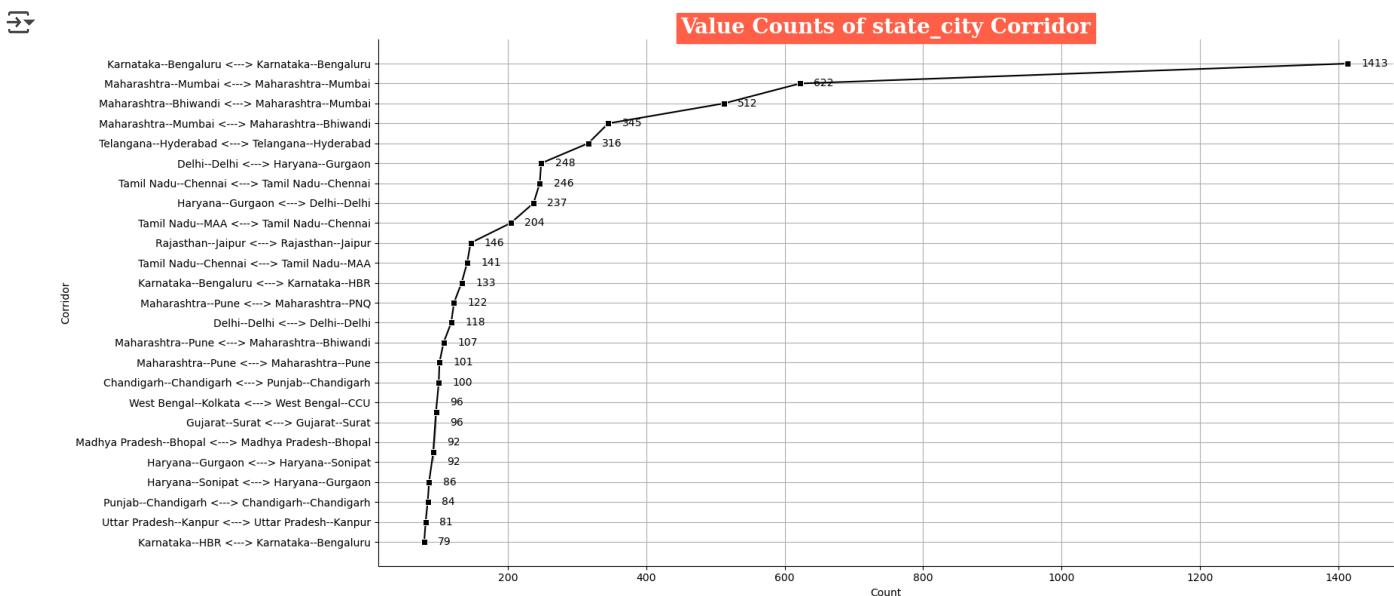
- The route between **Bangalore_Nelamangala_H** to **Bengaluru_KGAirport_HB,Bengaluru_Bomsndra_HB** sees the highest package volume, with 151 and 127 packages sent respectively.
- Bengaluru_Bommashandra_HB** to **Bengaluru_KGAirport_HB** is also popular, with 121 packages sent.
- Bengaluru_KGAirport_HB** to **Bangalore_Nelamangala_H** has moderate activity, with 108 packages sent.

4. The data indicates Bengaluru's importance as a transportation hub Corridor within **Karnataka**, handling significant package traffic.

```
1 de['state_corridor'] = de['source_state']+---+de['source_city'] +' <--> '+ de['destination_state']+---+de['destination_city']
2 de['state_corridor'].value_counts()
```

```
❖ state_corridor
Karnataka--Bengaluru <--> Karnataka--Bengaluru 1413
Maharashtra--Mumbai <--> Maharashtra--Mumbai 622
Maharashtra--Bhiwandi <--> Maharashtra--Mumbai 512
Maharashtra--Mumbai <--> Maharashtra--Bhiwandi 345
Telangana--Hyderabad <--> Telangana--Hyderabad 316
...
Gujarat--Jetpur <--> Gujarat--Dhoraji 1
Andhra Pradesh--Anakapalle <--> Andhra Pradesh--Visakhapatnam 1
Andhra Pradesh--Narsipatna <--> Andhra Pradesh--Anakapalle 1
West Bengal--MirzapurWB <--> West Bengal--Kolkata 1
Uttar Pradesh--Anandnagar <--> Uttar Pradesh--Gorakhpur 1
Name: count, Length: 2302, dtype: int64
```

```
1 state_corridor_counts = de['state_corridor'].value_counts()[:25]
2
3 plt.figure(figsize=(18,8))
4 sns.lineplot(y=state_corridor_counts.index, x=state_corridor_counts.values, marker='s', color='k')
5 plt.title('Value Counts of state_city Corridor', fontsize=20, fontfamily='serif', fontweight='bold', backgroundcolor='tomato', color='w')
6 plt.ylabel('Corridor')
7 plt.xlabel('Count')
8 plt.tight_layout()
9 sns.despine()
10 plt.grid(True)
11
12 for i, count in enumerate(state_corridor_counts.values):
13     plt.text(count+20, state_corridor_counts.index[i], str(count), ha='left', va='center')
14
15 plt.show()
```



```
1 de['city_corridor'] = de['source_city']+---+de['source_place'] +' <--> '+ de['destination_city']+---+de['destination_place']
2 display(de['city_corridor'].value_counts())
```

```

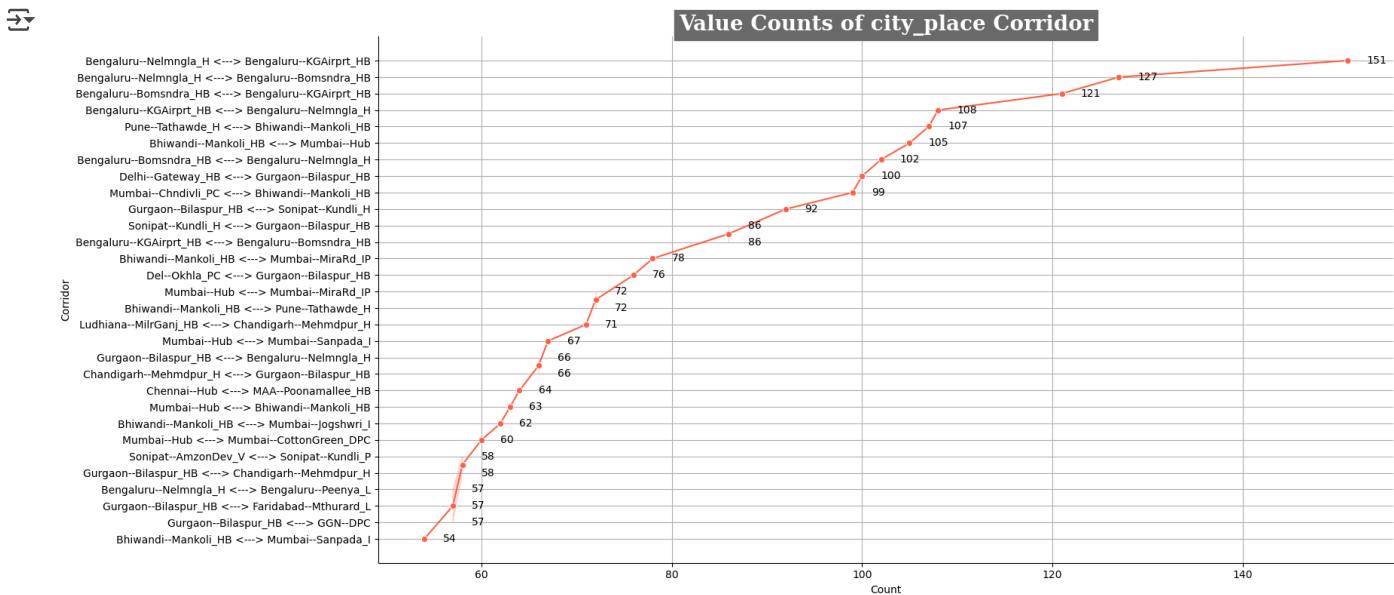
city_corridor
Bengaluru--Nelmngla_H <---> Bengaluru--KGAirprt_HB 151
Bengaluru--Nelmngla_H <---> Bengaluru--Bomsndra_HB 127
Bengaluru--Bomsndra_HB <---> Bengaluru--KGAirprt_HB 121
Bengaluru--KGAirprt_HB <---> Bengaluru--Nelmngla_H 108
Pune--Tathawde_H <---> Bhiwandi--Mankoli_HB 107
...
Ongole--SubhVRTL_I <---> Kandukur--LICOOffce_D 1
Madnapalle--PngnrRd_D <---> Palamaner--Lakshmi_D 1
Dharmavram--SaiNgr_D <---> Kadiri--GVManu_D 1
Baharampur--Chuanpur_I <---> Chapra--NagarDPP_D 1
Jaipur--NgrNigam_DC <---> Jaipur--Central_D_1 1
Name: count, Length: 2741, dtype: int64

```

```

1 city_corridor_counts = de['city_corridor'].value_counts()[:30]
2
3 plt.figure(figsize=(18,8))
4 sns.lineplot(y=city_corridor_counts.index, x=city_corridor_counts.values, marker='o', color='tomato')
5 plt.title('Value Counts of city_place Corridor', fontsize=20, fontfamily='serif', fontweight='bold', backgroundcolor='dimgray', color='w')
6 plt.ylabel('Corridor')
7 plt.xlabel('Count')
8 plt.tight_layout()
9 sns.despine()
10 plt.grid(True)
11
12 for i, count in enumerate(city_corridor_counts.values):
13     plt.text(count+2, city_corridor_counts.index[i], str(count), ha='left', va='center')
14
15 plt.show()

```



💡 Insights:

- Maharashtra, Karnataka, Haryana, and Tamil Nadu serve as key starting and ending locations for delivery services.
- Mumbai, Gurgaon, Delhi, and Bengaluru are major metropolitan centers from where many deliveries originate.
- A large proportion of nationwide deliveries are destined for Mumbai, Bengaluru, Gurgaon, and Delhi.

```

1 # 4. Extracting features like month, year, day, etc. from Trip_creation_time
2 de['trip_creation_month'] = de['trip_creation_time'].dt.month
3 de['trip_creation_year'] = de['trip_creation_time'].dt.year
4 de['trip_creation_day'] = de['trip_creation_time'].dt.day
5 de['trip_creation_hour'] = de['trip_creation_time'].dt.hour
6 de['trip_creation_weekday'] = de['trip_creation_time'].dt.weekday
7 de['trip_creation_week'] = de['trip_creation_time'].dt.isocalendar().week
8 de

```

| | segment_key | trip_uuid | data | route_type | trip_creation_time | source |
|-------|--|-------------------------|-------------------------|------------|--------------------|----------------------------|
| 0 | 153671041653548748+IND209304AAA+IND00000ACB | trip-153671041653548748 | trip-153671041653548748 | training | FTL | 2018-09-12 00:00:16.535741 |
| 1 | 153671041653548748+IND462022AAA+IND209304AAA | trip-153671041653548748 | trip-153671041653548748 | training | FTL | 2018-09-12 00:00:16.535741 |
| 2 | 153671042288605164+IND561203AAB+IND562101AAA | trip-153671042288605164 | trip-153671042288605164 | training | Carting | 2018-09-12 00:00:22.886430 |
| 3 | 153671042288605164+IND572101AAA+IND561203AAB | trip-153671042288605164 | trip-153671042288605164 | training | Carting | 2018-09-12 00:00:22.886430 |
| 4 | 153671043369099517+IND00000ACB+IND160002AAC | trip-153671043369099517 | trip-153671043369099517 | training | FTL | 2018-09-12 00:00:33.691250 |
| ... | ... | ... | ... | ... | ... | ... |
| 26217 | 153861115439069069+IND628204AAA+IND627657AAA | trip-153861115439069069 | trip-153861115439069069 | test | Carting | 2018-10-03 23:59:14.390954 |
| 26218 | 153861115439069069+IND628613AAA+IND627005AAA | trip-153861115439069069 | trip-153861115439069069 | test | Carting | 2018-10-03 23:59:14.390954 |
| 26219 | 153861115439069069+IND628801AAA+IND628204AAA | trip-153861115439069069 | trip-153861115439069069 | test | Carting | 2018-10-03 23:59:14.390954 |
| 26220 | 153861118270144424+IND583119AAA+IND583101AAA | trip-153861118270144424 | trip-153861118270144424 | test | FTL | 2018-10-03 23:59:42.701692 |
| 26221 | 153861118270144424+IND583201AAA+IND583119AAA | trip-153861118270144424 | trip-153861118270144424 | test | FTL | 2018-10-03 23:59:42.701692 |

26222 rows × 37 columns



❖ In-Depth Analysis

```
1 new_df = de.copy()
```

```
1 new_df.columns
```

```
2 Index(['segment_key', 'trip_uuid', 'data', 'route_type', 'trip_creation_time',
       'source_name', 'destination_name', 'od_start_time', 'od_end_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', 'segment_actual_time_sum',
       'segment_osrm_time_sum', 'segment_osrm_distance_sum', 'od_total_time',
       'od_time_diff_hour', 'source_city', 'source_place', 'source_state',
       'destination_city', 'destination_place', 'destination_state',
       'corridor', 'state_corridor', 'city_corridor', 'trip_creation_month',
       'trip_creation_year', 'trip_creation_day', 'trip_creation_hour',
       'trip_creation_weekday', 'trip_creation_week'],
      dtype='object')
```

```
1 new_df.sample(2)
```

| | segment_key | trip_uuid | data | route_type | trip_creation_time | source_na |
|------|--|-------------------------|----------|------------|----------------------------|-----------------------------|
| 3330 | 153694740553026744+IND574104AAA+IND574216AAA | trip-153694740553026744 | training | Carting | 2018-09-14 17:50:05.530477 | Karkala_MarketRd (Karnatal |
| 3755 | 153696604384005633+IND507117AAB+IND507303AAA | trip-153696604384005633 | training | FTL | 2018-09-14 23:00:43.840397 | Manuguru_AskNagar (Telangar |

```

1 create_trip_dict={
2   'data' : 'first',
3   'route_type' : 'first',
4   'od_start_time':'first',
5   'od_end_time':'last',
6   'od_time_diff_hour' : 'sum',
7   'trip_creation_time' : 'first',
8   'trip_creation_month' : 'first',
9   'trip_creation_year' : 'first',
10  'trip_creation_day' : 'first',
11  'trip_creation_hour' : 'first',
12  'trip_creation_weekday' : 'first',
13  'trip_creation_week' : 'first',
14  'start_scan_to_end_scan' : 'sum',
15  'actual_distance_to_destination' : 'sum',
16  'actual_time' : 'sum',
17  'osrm_time' : 'sum',
18  'osrm_distance' : 'sum',
19  'segment_actual_time': 'sum',
20  'segment_osrm_time': 'sum',
21  'segment_osrm_distance': 'sum',
22  'segment_actual_time_sum': 'sum',
23  'segment_osrm_time_sum': 'sum',
24  'segment_osrm_distance_sum': 'sum',
25  'source_name': 'first',
26  'source_city':'first',
27  'source_state':'first',
28  'source_place':'first',
29  'destination_name': 'first',
30  'destination_city':'first',
31  'destination_state':'first',
32  'destination_place':'first',
33  'corridor':'first',
34  'state_corridor':'first',
35  'city_corridor':'first'
36 }
37
38 trip_agg_df = new_df.groupby('trip_uuid').agg(create_trip_dict).reset_index()
39 trip_agg_df

```

| | trip_uuid | data | route_type | od_start_time | od_end_time | od_time_diff_hour | trip_creation_time | trip_creation_i |
|-------|--------------------|---------------|------------|----------------------------|----------------------------|-------------------|----------------------------|-----------------|
| 0 | 153671041653548748 | trip-training | FTL | 2018-09-12 16:39:46.858469 | 2018-09-12 16:39:46.858469 | 37.668497 | 2018-09-12 00:00:16.535741 | |
| 1 | 153671042288605164 | trip-training | Carting | 2018-09-12 02:03:09.655591 | 2018-09-12 02:03:09.655591 | 3.026865 | 2018-09-12 00:00:22.886430 | |
| 2 | 153671043369099517 | trip-training | FTL | 2018-09-14 03:40:17.106733 | 2018-09-14 03:40:17.106733 | 65.572709 | 2018-09-12 00:00:33.691250 | |
| 3 | 153671046011330457 | trip-training | Carting | 2018-09-12 00:01:00.113710 | 2018-09-12 01:41:29.809822 | 1.674916 | 2018-09-12 00:01:00.113710 | |
| 4 | 153671052974046625 | trip-training | FTL | 2018-09-12 00:02:09.740725 | 2018-09-12 03:54:43.114421 | 11.972484 | 2018-09-12 00:02:09.740725 | |
| ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 14782 | 153861095625827784 | trip-test | Carting | 2018-10-03 23:55:56.258533 | 2018-10-04 06:41:25.409035 | 4.300482 | 2018-10-03 23:55:56.258533 | |
| 14783 | 153861104386292051 | trip-test | Carting | 2018-10-03 23:57:23.863155 | 2018-10-04 00:57:59.294434 | 1.009842 | 2018-10-03 23:57:23.863155 | |
| 14784 | 153861106442901555 | trip-test | Carting | 2018-10-04 02:51:27.075797 | 2018-10-04 02:51:27.075797 | 7.035331 | 2018-10-03 23:57:44.429324 | |
| 14785 | 153861115439069069 | trip-test | Carting | 2018-10-03 23:59:14.390954 | 2018-10-04 02:29:04.272194 | 5.808548 | 2018-10-03 23:59:14.390954 | |
| 14786 | 153861118270144424 | trip-test | FTL | 2018-10-04 03:58:40.726547 | 2018-10-04 03:58:40.726547 | 5.906793 | 2018-10-03 23:59:42.701692 | |

14787 rows × 35 columns

```
1 numerical_columns = trip_agg_df.select_dtypes(include=[np.float32, np.float64])
2 numerical_columns
```

| | od_time_diff_hour | start_scan_to_end_scan | actual_distance_to_destination | actual_time | osrm_time | osrm_distance | segment_act |
|-------|-------------------|------------------------|--------------------------------|-------------|-----------|---------------|-------------|
| 0 | 37.668497 | 2259.0 | 824.732849 | 1562.0 | 717.0 | 991.352295 | |
| 1 | 3.026865 | 180.0 | 73.186905 | 143.0 | 68.0 | 85.111000 | |
| 2 | 65.572709 | 3933.0 | 1927.404297 | 3347.0 | 1740.0 | 2354.066650 | |
| 3 | 1.674916 | 100.0 | 17.175274 | 59.0 | 15.0 | 19.680000 | |
| 4 | 11.972484 | 717.0 | 127.448502 | 341.0 | 117.0 | 146.791794 | |
| ... | ... | ... | ... | ... | ... | ... | ... |
| 14782 | 4.300482 | 257.0 | 57.762333 | 83.0 | 62.0 | 73.462997 | |
| 14783 | 1.009842 | 60.0 | 15.513784 | 21.0 | 12.0 | 16.088200 | |
| 14784 | 7.035331 | 421.0 | 38.684837 | 282.0 | 48.0 | 58.903702 | |
| 14785 | 5.808548 | 347.0 | 134.723831 | 264.0 | 179.0 | 171.110306 | |
| 14786 | 5.906793 | 353.0 | 66.081528 | 275.0 | 68.0 | 80.578705 | |

14787 rows × 12 columns

```
1 numerical_columns.describe().T
```

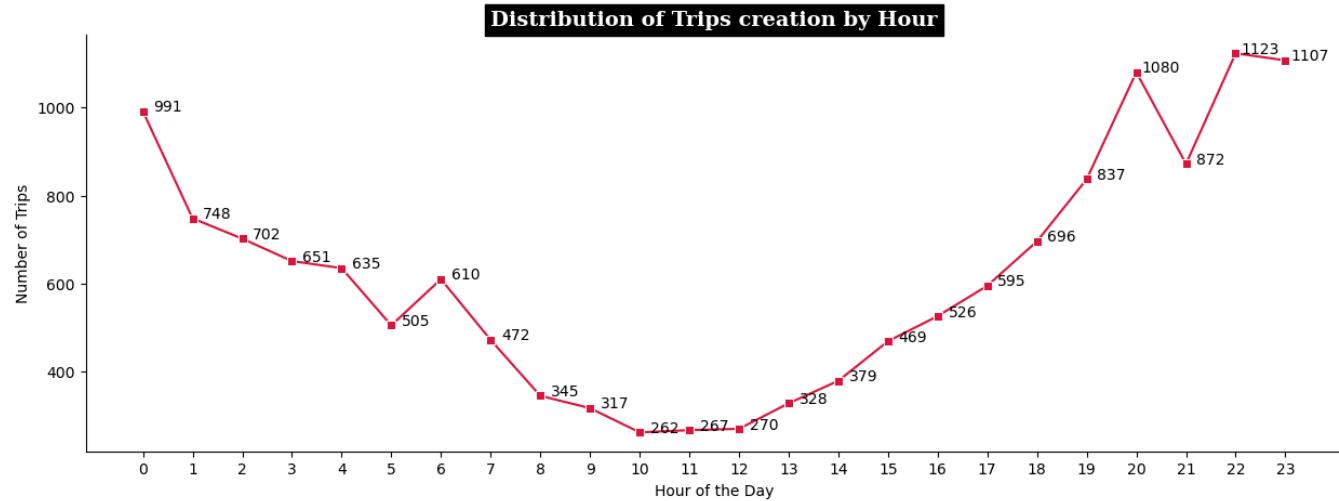
| | count | mean | std | min | 25% | 50% | 75% | max |
|--------------------------------|---------|------------|------------|-----------|------------|------------|------------|-------------|
| od_time_diff_hour | 14787.0 | 8.840187 | 10.978880 | 0.391024 | 2.494975 | 4.661846 | 10.558962 | 131.642533 |
| start_scan_to_end_scan | 14787.0 | 529.429016 | 658.254395 | 23.000000 | 149.000000 | 279.000000 | 632.000000 | 7898.000000 |
| actual_distance_to_destination | 14787.0 | 164.090195 | 305.502808 | 9.002461 | 22.777099 | 48.287895 | 163.591255 | 2186.531738 |
| actual_time | 14787.0 | 356.306000 | 561.517761 | 9.000000 | 67.000000 | 148.000000 | 367.000000 | 6265.000000 |
| osrm_time | 14787.0 | 160.990936 | 271.459229 | 6.000000 | 29.000000 | 60.000000 | 168.000000 | 2032.000000 |
| osrm_distance | 14787.0 | 203.887405 | 370.565460 | 9.072900 | 30.756900 | 65.302795 | 206.644203 | 2840.081055 |
| segment_actual_time | 14787.0 | 353.059174 | 556.364441 | 9.000000 | 66.000000 | 147.000000 | 364.000000 | 6230.000000 |
| segment_osrm_time | 14787.0 | 180.511597 | 314.678741 | 6.000000 | 30.000000 | 65.000000 | 184.000000 | 2564.000000 |
| segment_osrm_distance | 14787.0 | 222.705444 | 416.845642 | 9.072900 | 32.578850 | 69.784203 | 216.560608 | 3523.632324 |
| segment_actual_time_sum | 14787.0 | 353.059174 | 556.364441 | 9.000000 | 66.000000 | 147.000000 | 364.000000 | 6230.000000 |
| segment_osrm_time_sum | 14787.0 | 180.511597 | 314.678741 | 6.000000 | 30.000000 | 65.000000 | 184.000000 | 2564.000000 |
| segment_osrm_distance_sum | 14787.0 | 222.705444 | 416.845642 | 9.072900 | 32.578850 | 69.784203 | 216.560608 | 3523.632324 |

```
1 trip_agg_df.describe(include = object).T
```

| | count | unique | top | freq |
|-------------------|-------|--------|---|------|
| trip_uuid | 14787 | 14787 | trip-153671041653548748 | 1 |
| source_name | 14787 | 930 | Gurgaon_Bilaspur_HB (Haryana) | 1052 |
| source_city | 14787 | 713 | Bengaluru | 1700 |
| source_state | 14787 | 29 | Maharashtra | 2714 |
| source_place | 14787 | 788 | Bilaspur_HB | 1052 |
| destination_name | 14787 | 1042 | Gurgaon_Bilaspur_HB (Haryana) | 745 |
| destination_city | 14787 | 851 | Bengaluru | 1633 |
| destination_state | 14787 | 32 | Maharashtra | 2569 |
| destination_place | 14787 | 866 | Bilaspur_HB | 745 |
| corridor | 14787 | 1737 | Bangalore_Nelmngla_H (Karnataka) <--> Bengaluru_H | 151 |
| state_corridor | 14787 | 1366 | Karnataka--Bengaluru <--> Karnataka--Bengaluru | 1333 |
| city_corridor | 14787 | 1737 | Bengaluru--Nelmngla_H <--> Bengaluru--KGAirport | 151 |

```
1 trip_df = trip_agg_df.copy()
```

```
1 trip_creation_by_hour = trip_df.groupby(by='trip_creation_hour')['trip_uuid'].count().reset_index()
2
3 plt.figure(figsize=(15,5))
4 sns.lineplot(data=trip_creation_by_hour, x='trip_creation_hour', y='trip_uuid', marker='s', color='crimson')
5 plt.xticks(np.arange(0, 24))
6
7 for i, count in enumerate(trip_creation_by_hour['trip_uuid']):
8     plt.text(trip_creation_by_hour['trip_creation_hour'][i]+0.5, count, count, ha='center')
9
10 plt.title('Distribution of Trips creation by Hour', fontsize=14, fontfamily='serif', fontweight='bold', backgroundcolor='k', color='w')
11 plt.xlabel('Hour of the Day')
12 plt.ylabel('Number of Trips')
13 sns.despine()
14 plt.show()
```



```
1 trip_df.sample()
```

```
trip_uuid      data route_type od_start_time      od_end_time od_time_diff_hour trip_creation_time trip_creation_m
7917 153765583943264913 trip-training Carting 2018-09-22 22:37:19.432896 23:19:42.965506 0.706537 2018-09-22 22:37:19.432896
```

```
1 trip_df.trip_creation_year.value_counts()
```

```
trip_creation_year
2018    14787
Name: count, dtype: int64
```

```
1 trip_df.trip_creation_month.value_counts()
```

```
trip_creation_month
9       13011
10      1776
Name: count, dtype: int64
```

```
1 trip_df['trip_creation_month'].value_counts(normalize = True) * 100
```

```
trip_creation_month
9      87.98945
10     12.01055
Name: proportion, dtype: float64
```

```
1 trip_df.trip_creation_week.value_counts()
```

```
trip_creation_week
38      5001
39      4402
37      3608
40      1776
Name: count, dtype: Int64
```

```
1 trip_df.trip_creation_weekday.value_counts(ascending=True)
```

```
trip_creation_weekday
6      1753
0      1980
1      2035
4      2057
3      2103
5      2128
2      2731
Name: count, dtype: int64
```

```
1 trip_df['trip_creation_day_week'] = trip_df['trip_creation_time'].dt.day_name()
```

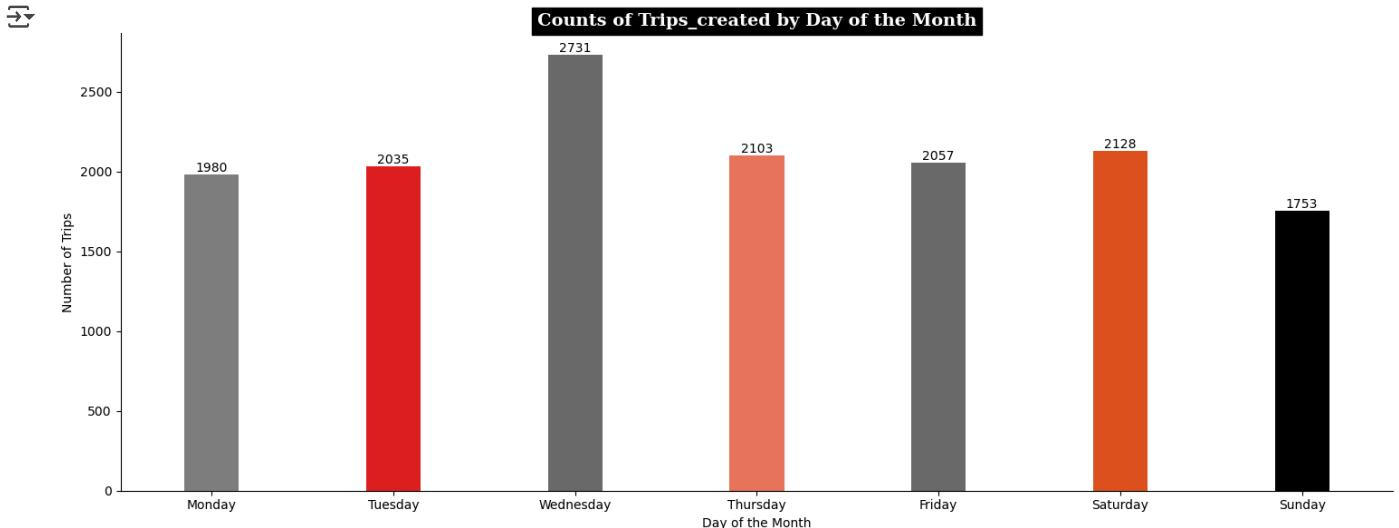
```
1 trip_df.trip_creation_day.value_counts()
```

```
trip_creation_day
18    791
15    783
13    750
12    747
21    740
22    740
17    722
14    712
20    703
25    695
26    683
19    674
24    658
27    650
23    631
3     627
16    616
28    605
29    605
1     600
2     549
30    506
Name: count, dtype: int64
```

```
1 trip_df.sample()
```

| | trip_uuid | data | route_type | od_start_time | od_end_time | od_time_diff_hour | trip_creation_time | trip_creation_m |
|------|--------------------|----------------|------------|----------------------------|----------------------------|-------------------|----------------------------|-----------------|
| 8403 | 153772987770593762 | trip- training | FTL | 2018-09-24 00:14:22.226774 | 2018-09-24 02:42:00.976996 | 3.847563 | 2018-09-23 19:11:17.706293 | |

```
1 plt.figure(figsize=(15,6))
2
3 weekday_order = ['Monday', 'Tuesday', 'Wednesday', 'Thursday', 'Friday', 'Saturday', 'Sunday']
4 day_counts = trip_df['trip_creation_day_week'].value_counts().reindex(weekday_order)
5
6 sns.barplot(x=day_counts.index, y=day_counts.values, palette=cp, width=0.3)
7 for i, count in enumerate(day_counts.values):
8     plt.text(i, count, str(count), ha='center', va='bottom')
9 plt.title('Counts of Trips_created by Day of the Month', fontsize=14, fontfamily='serif', fontweight='bold', backgroundcolor='k', color='w')
10 plt.xlabel('Day of the Month')
11 plt.ylabel('Number of Trips')
12 plt.tight_layout()
13 sns.despine()
14 plt.show()
```

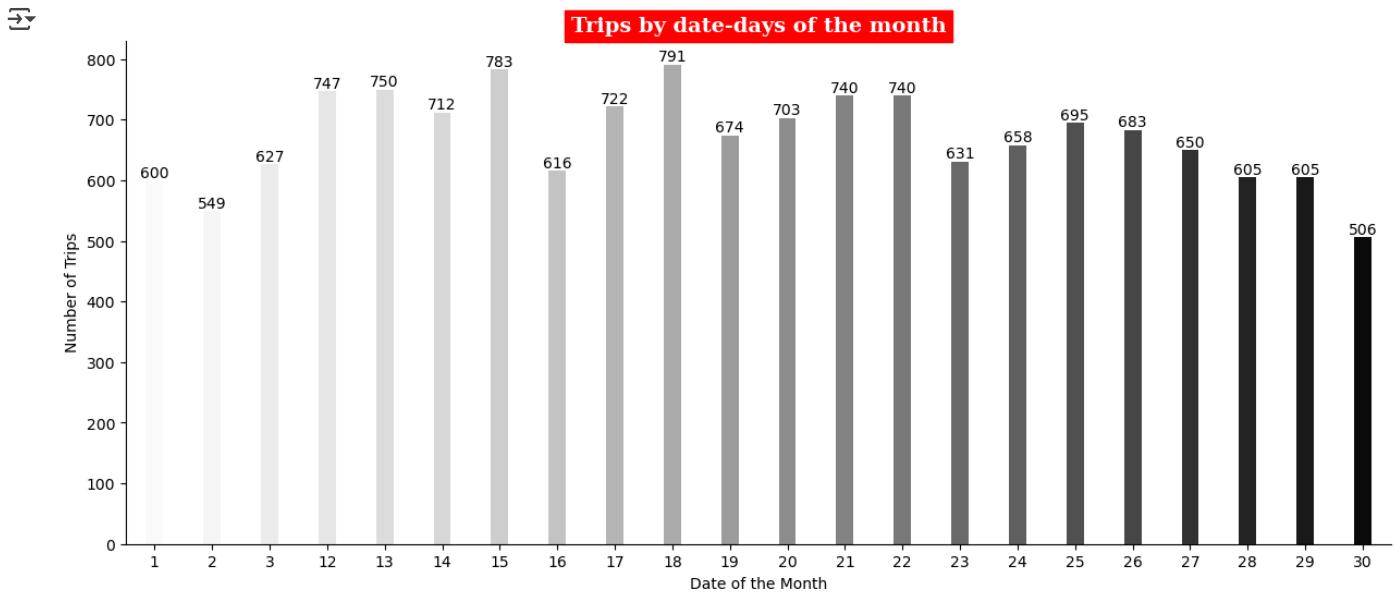


```

1 trip_df['trip_creation_dayofdate'] = trip_df['trip_creation_time'].dt.day

1 trips_by_dateday = trip_df.groupby(by = 'trip_creation_dayofdate')['trip_uuid'].count().to_frame().reset_index()
2
3 plt.figure(figsize = (15, 6))
4 sns.barplot(data = trip_df,x = trips_by_dateday['trip_creation_dayofdate'],y = trips_by_dateday['trip_uuid'], palette='Greys',width=0
5 for i, count in enumerate(trips_by_dateday['trip_uuid']):
6     plt.text(i, count, str(count), ha='center', va='bottom')
7 plt.title('Trips by date-days of the month',fontsize=14,fontfamily='serif',fontweight='bold',backgroundcolor='r',color='w')
8 plt.xlabel('Date of the Month')
9 plt.ylabel('Number of Trips')
10 sns.despine()
11 plt.show()

```



Outlier treatment

```
1 numerical_columns
```

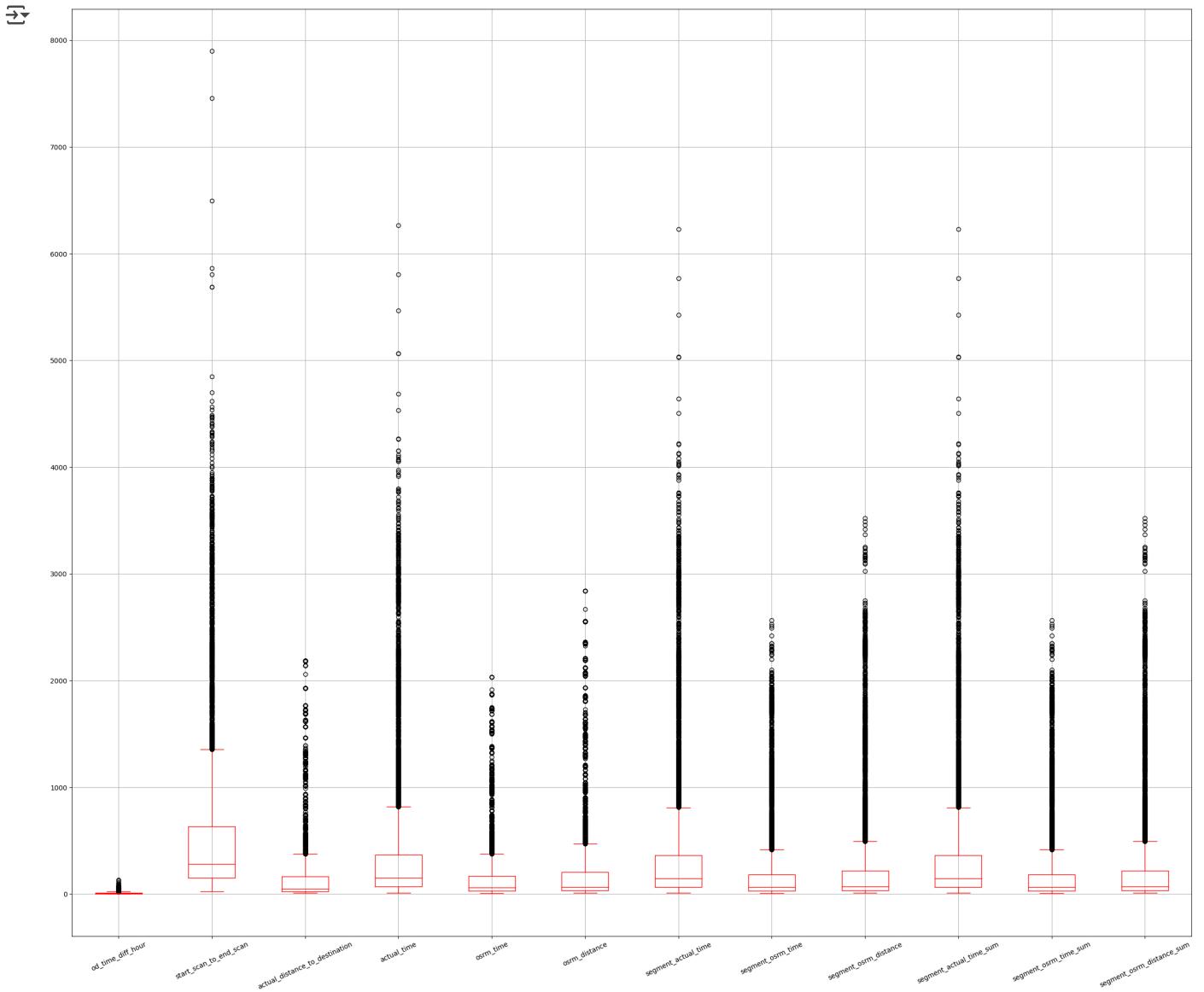
| | od_time_diff_hour | start_scan_to_end_scan | actual_distance_to_destination | actual_time | osrm_time | osrm_distance | segment_act |
|-------|-------------------|------------------------|--------------------------------|-------------|-----------|---------------|-------------|
| 0 | 37.668497 | 2259.0 | 824.732849 | 1562.0 | 717.0 | 991.352295 | |
| 1 | 3.026865 | 180.0 | 73.186905 | 143.0 | 68.0 | 85.111000 | |
| 2 | 65.572709 | 3933.0 | 1927.404297 | 3347.0 | 1740.0 | 2354.066650 | |
| 3 | 1.674916 | 100.0 | 17.175274 | 59.0 | 15.0 | 19.680000 | |
| 4 | 11.972484 | 717.0 | 127.448502 | 341.0 | 117.0 | 146.791794 | |
| ... | ... | ... | ... | ... | ... | ... | ... |
| 14782 | 4.300482 | 257.0 | 57.762333 | 83.0 | 62.0 | 73.462997 | |
| 14783 | 1.009842 | 60.0 | 15.513784 | 21.0 | 12.0 | 16.088200 | |
| 14784 | 7.035331 | 421.0 | 38.684837 | 282.0 | 48.0 | 58.903702 | |
| 14785 | 5.808548 | 347.0 | 134.723831 | 264.0 | 179.0 | 171.110306 | |
| 14786 | 5.906793 | 353.0 | 66.081528 | 275.0 | 68.0 | 80.578705 | |

14787 rows × 12 columns

```

1 plt.figure(figsize=(30, 25))
2 numerical_columns.boxplot(rot=25, figsize=(35,20), color = 'r')
3 plt.grid('off')
4 plt.show()

```



```
1 numerical_columns.columns
```

```
2 [Index(['od_time_diff_hour', 'start_scan_to_end_scan',
       'actual_distance_to_destination', 'actual_time', 'osrm_time',
       'osrm_distance', 'segment_actual_time', 'segment_osrm_time',
       'segment_osrm_distance', 'segment_actual_time_sum',
       'segment_osrm_time_sum', 'segment_osrm_distance_sum'],
      dtype='object')]
```

```
1 num_cols = numerical_columns.columns.tolist()
2 num_cols
```

```
3 ['od_time_diff_hour',
   'start_scan_to_end_scan',
   'actual_distance_to_destination',
   'actual_time',
   'osrm_time',
```

```
'osrm_distance',
'segment_actual_time',
'segment_osrm_time',
'segment_osrm_distance',
'segment_actual_time_sum',
'segment_osrm_time_sum',
'segment_osrm_distance_sum']
```

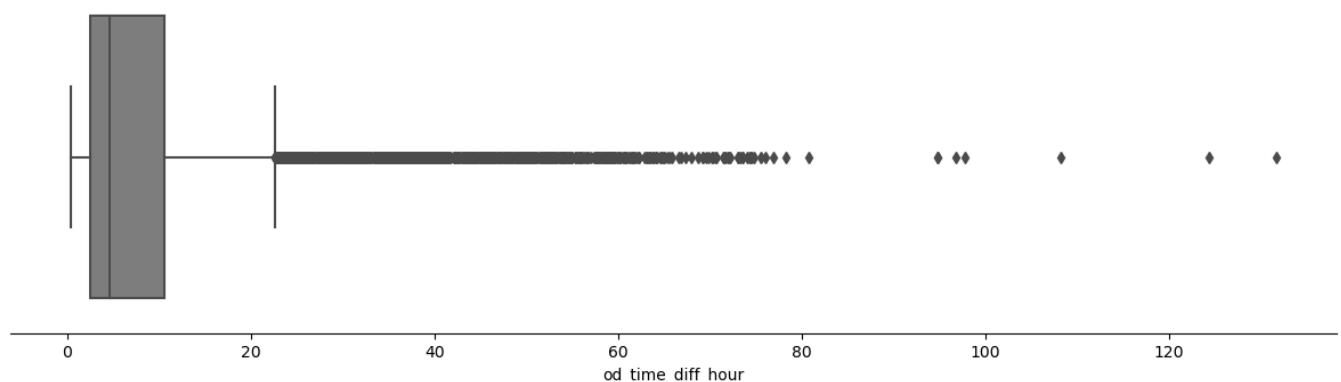
```
1 # obtain the first quartile
2 Q1 = numerical_columns.quantile(0.25)
3
4 # obtain the third quartile
5 Q3 = numerical_columns.quantile(0.75)
6
7 # obtain the IQR
8 IQR = Q3 - Q1
9
10 # print the IQR
11 print(IQR)
```

```
→ od_time_diff_hour          8.063987
start_scan_to_end_scan      483.000000
actual_distance_to_destination 140.814157
actual_time                  300.000000
osrm_time                    139.000000
osrm_distance                175.887303
segment_actual_time          298.000000
segment_osrm_time             154.000000
segment_osrm_distance         183.981758
segment_actual_time_sum       298.000000
segment_osrm_time_sum         154.000000
segment_osrm_distance_sum     183.981758
dtype: float64
```

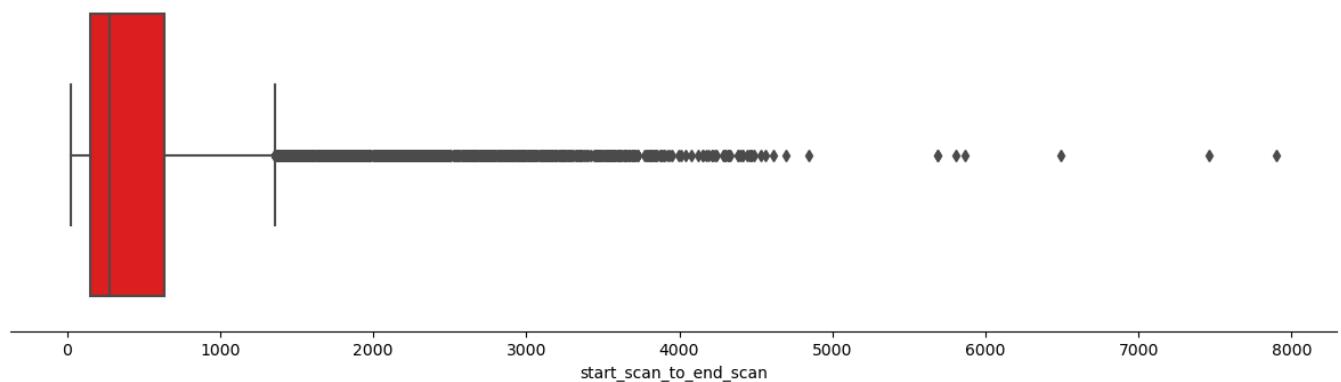
```
1 for i,col in enumerate(numerical_columns):
2     plt.figure(figsize=(15,4))
3     sns.boxplot(x=col, data=numerical_columns,color=cp[i])
4     sns.despine(left=True)
5     plt.yticks([])
6     plt.title(f'Boxplot of {col}',fontfamily='serif',fontweight='bold',fontsize=12,backgroundcolor=cp[i],color='w')
7     plt.show()
```



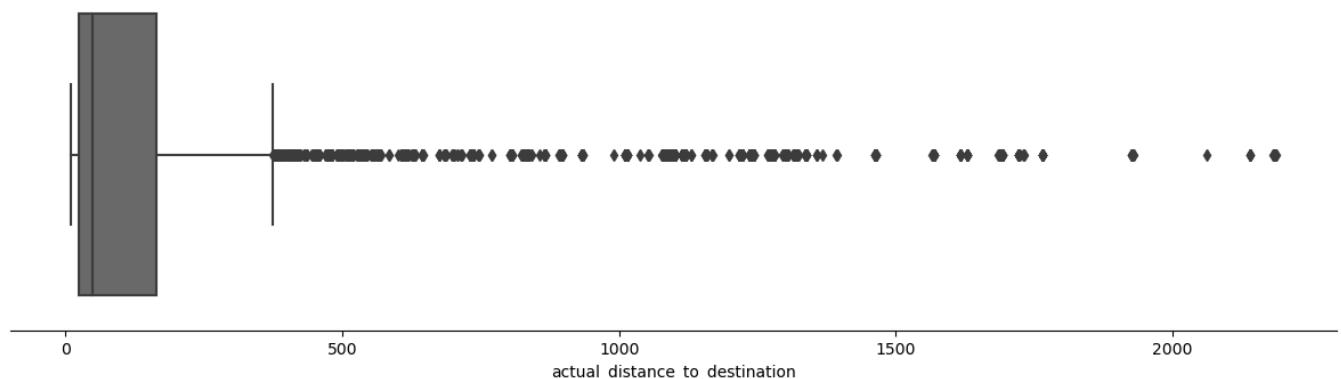
Boxplot of od_time_diff_hour



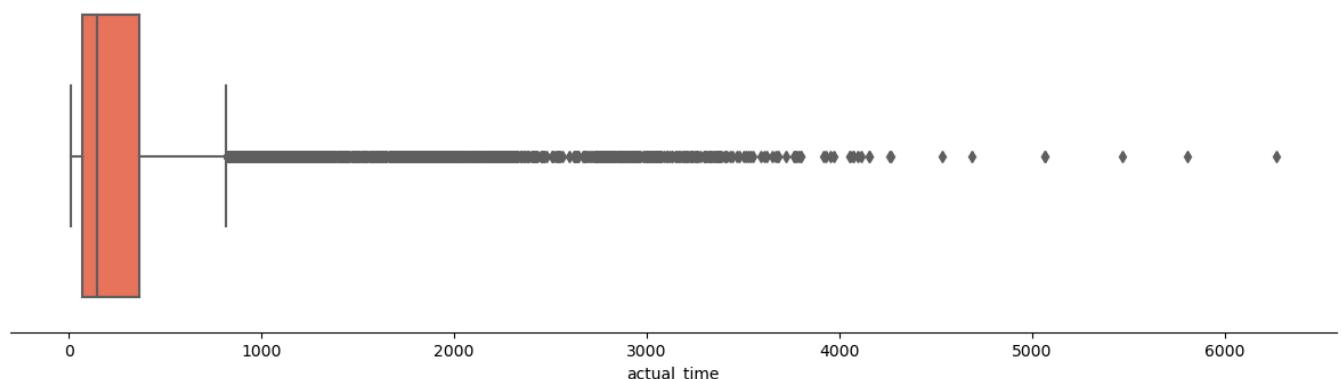
Boxplot of start_scan_to_end_scan



Boxplot of actual_distance_to_destination

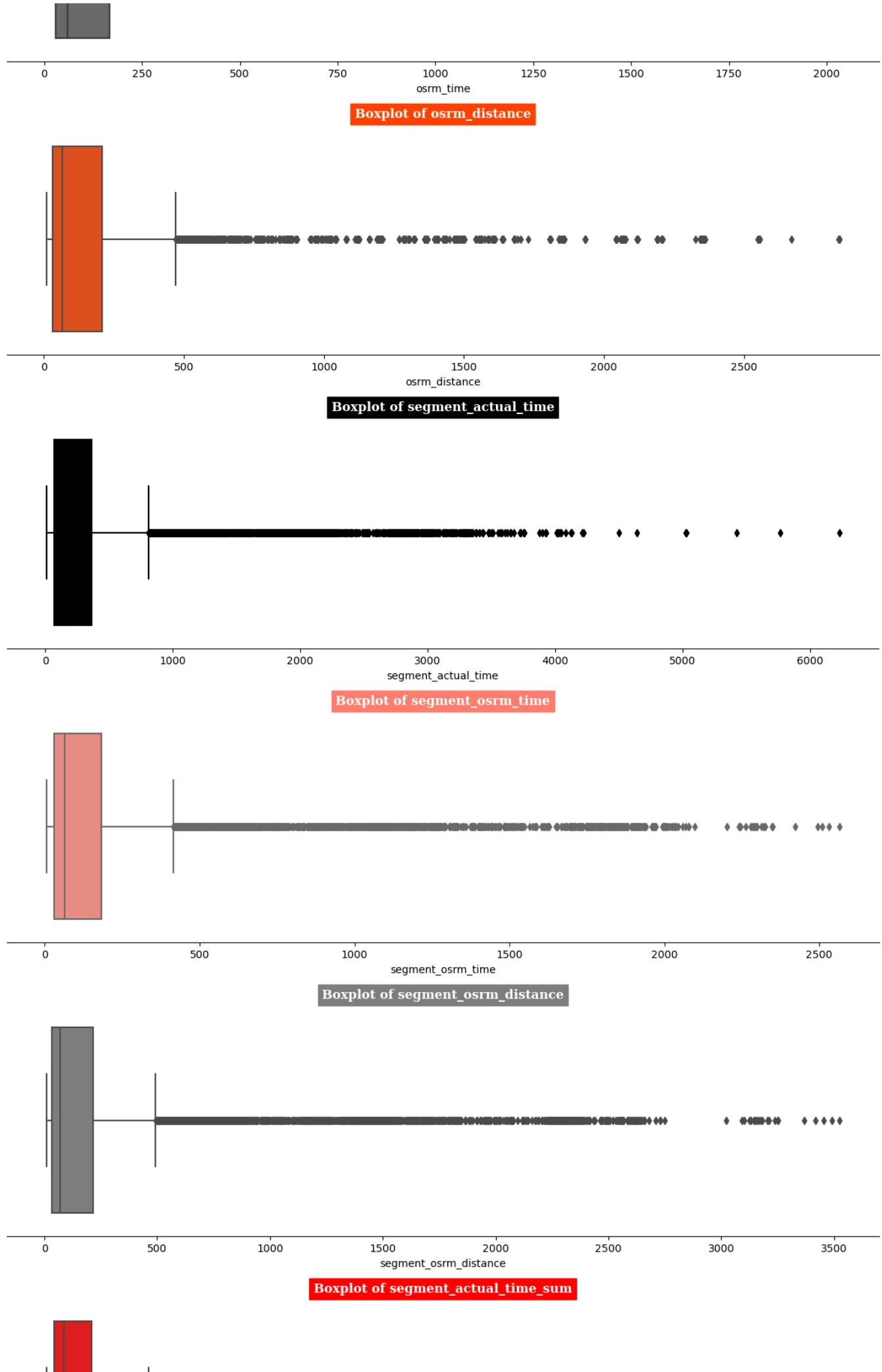


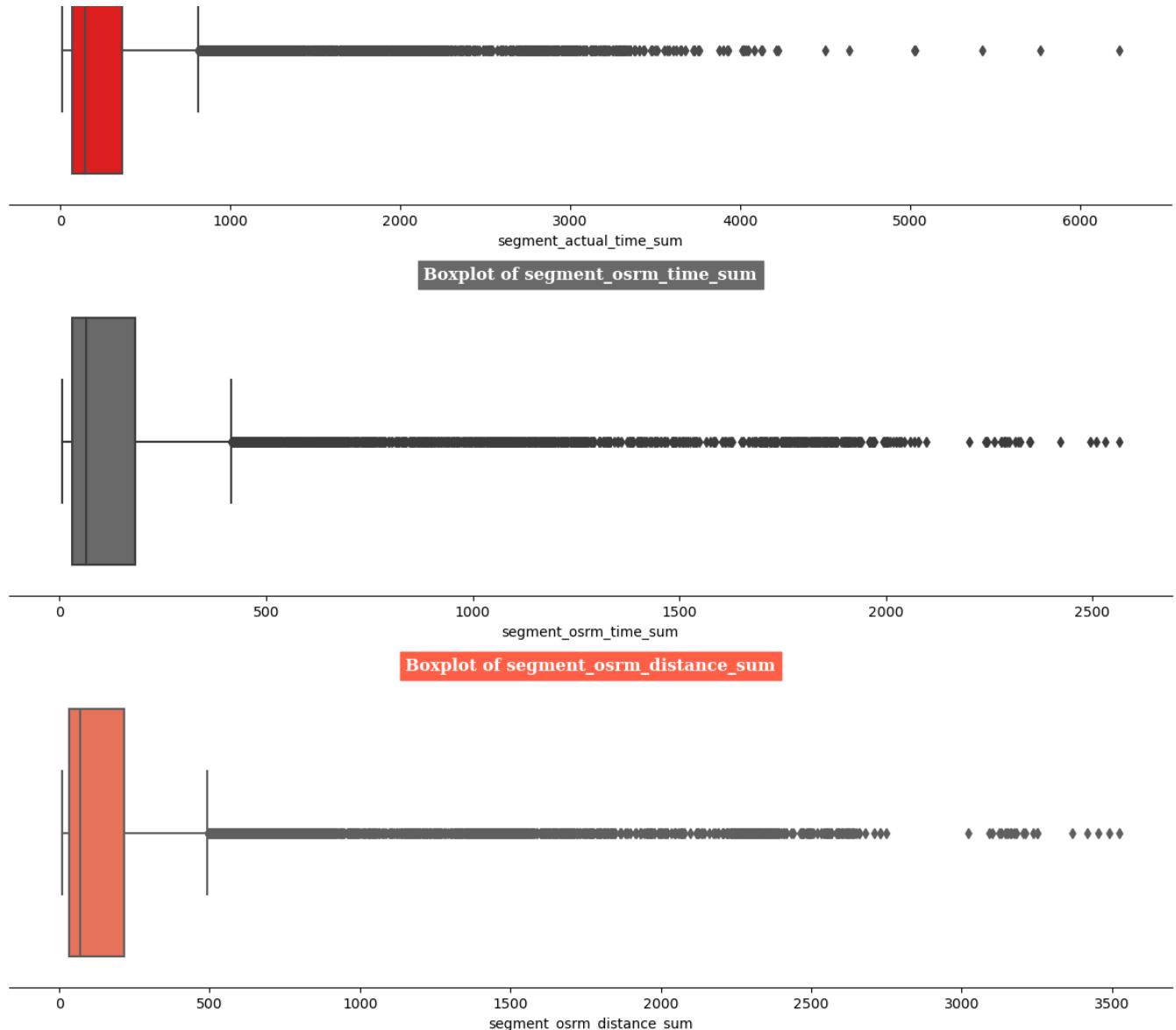
Boxplot of actual_time



Boxplot of osrm_time







Outlier Removal

```
1 for i, col in enumerate(numerical_columns):
2
3     data = trip_df[col]
4     display(data.to_frame())
5
6
7     Q1 = np.percentile(data, 25)
8     Q3 = np.percentile(data, 75)
9     IQR = Q3 - Q1
10
11
12     lower_bound = Q1 - (1.5 * IQR)
13     upper_bound = Q3 + (1.5 * IQR)
14
15     clipped_data = np.clip(data, lower_bound, upper_bound)
16     print(f'Clipped data of {col}')
17     display(clipped_data.to_frame())
18     print()
19
20     # Plot boxplot of the clipped data
21     plt.figure(figsize=(15, 4))
22     plt.subplot(121)
23     sns.boxplot(x=clipped_data, color=cp[i])
24     sns.despine(left=True)
25     plt.yticks([])
26     plt.title(f'Boxplot of clipped {col}', fontfamily='serif', fontweight='bold', fontsize=12, backgroundcolor=cp[i], color='w')
27
28     filtered_data = data.loc[(data >= lower_bound) | (data <= upper_bound)]
29     print(f'Filtered data of {col}')
30     display(filtered_data.to_frame())
31     print()
32
33     plt.subplot(122)
34     sns.boxplot(x=filtered_data, color=cp[i])
35     sns.despine(left=True)
36     plt.yticks([])
37     plt.title(f'Boxplot of filtered {col}', fontfamily='serif', fontweight='bold', fontsize=12, backgroundcolor=cp[i], color='w')
38
39     plt.show()
```

| | od_time_diff_hour |
|--------------|-------------------|
| 0 | 37.668497 |
| 1 | 3.026865 |
| 2 | 65.572709 |
| 3 | 1.674916 |
| 4 | 11.972484 |
| ... | ... |
| 14782 | 4.300482 |
| 14783 | 1.009842 |
| 14784 | 7.035331 |
| 14785 | 5.808548 |
| 14786 | 5.906793 |

14787 rows × 1 columns

Clipped data of od_time_diff_hour

| | od_time_diff_hour |
|--------------|-------------------|
| 0 | 22.654942 |
| 1 | 3.026865 |
| 2 | 22.654942 |
| 3 | 1.674916 |
| 4 | 11.972484 |
| ... | ... |
| 14782 | 4.300482 |
| 14783 | 1.009842 |
| 14784 | 7.035331 |
| 14785 | 5.808548 |
| 14786 | 5.906793 |

14787 rows × 1 columns

Filtered data of od_time_diff_hour

| | od_time_diff_hour |
|--------------|-------------------|
| 0 | 37.668497 |
| 1 | 3.026865 |
| 2 | 65.572709 |
| 3 | 1.674916 |
| 4 | 11.972484 |
| ... | ... |
| 14782 | 4.300482 |
| 14783 | 1.009842 |
| 14784 | 7.035331 |
| 14785 | 5.808548 |
| 14786 | 5.906793 |

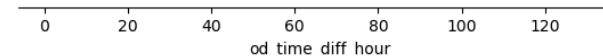
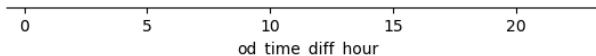
14787 rows × 1 columns

Boxplot of clipped od_time_diff_hour



Boxplot of filtered od_time_diff_hour



**start_scan_to_end_scan**

| | |
|-------|--------|
| 0 | 2259.0 |
| 1 | 180.0 |
| 2 | 3933.0 |
| 3 | 100.0 |
| 4 | 717.0 |
| ... | ... |
| 14782 | 257.0 |
| 14783 | 60.0 |
| 14784 | 421.0 |
| 14785 | 347.0 |
| 14786 | 353.0 |

14787 rows × 1 columns

Clipped data of start_scan_to_end_scan

start_scan_to_end_scan

| | |
|-------|--------|
| 0 | 1356.5 |
| 1 | 180.0 |
| 2 | 1356.5 |
| 3 | 100.0 |
| 4 | 717.0 |
| ... | ... |
| 14782 | 257.0 |
| 14783 | 60.0 |
| 14784 | 421.0 |
| 14785 | 347.0 |
| 14786 | 353.0 |

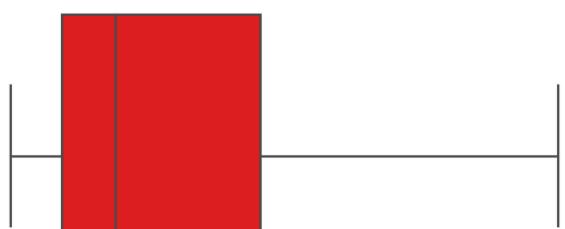
14787 rows × 1 columns

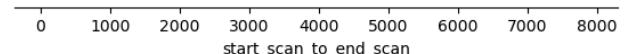
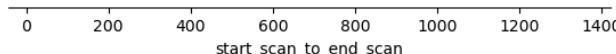
Filtered data of start_scan_to_end_scan

start_scan_to_end_scan

| | |
|-------|--------|
| 0 | 2259.0 |
| 1 | 180.0 |
| 2 | 3933.0 |
| 3 | 100.0 |
| 4 | 717.0 |
| ... | ... |
| 14782 | 257.0 |
| 14783 | 60.0 |
| 14784 | 421.0 |
| 14785 | 347.0 |
| 14786 | 353.0 |

14787 rows × 1 columns

Boxplot of clipped start_scan_to_end_scan**Boxplot of filtered start_scan_to_end_scan**

**actual_distance_to_destination**

| | |
|--------------|-------------|
| 0 | 824.732849 |
| 1 | 73.186905 |
| 2 | 1927.404297 |
| 3 | 17.175274 |
| 4 | 127.448502 |
| ... | ... |
| 14782 | 57.762333 |
| 14783 | 15.513784 |
| 14784 | 38.684837 |
| 14785 | 134.723831 |
| 14786 | 66.081528 |

14787 rows × 1 columns

Clipped data of actual_distance_to_destination

actual_distance_to_destination

| | |
|--------------|------------|
| 0 | 374.812490 |
| 1 | 73.186905 |
| 2 | 374.812490 |
| 3 | 17.175274 |
| 4 | 127.448502 |
| ... | ... |
| 14782 | 57.762333 |
| 14783 | 15.513784 |
| 14784 | 38.684837 |
| 14785 | 134.723831 |
| 14786 | 66.081528 |

14787 rows × 1 columns

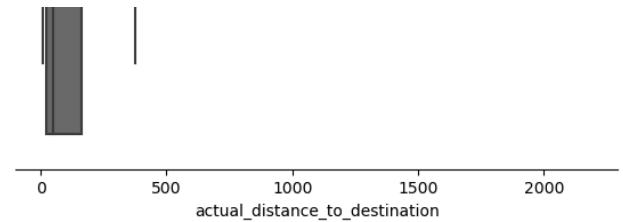
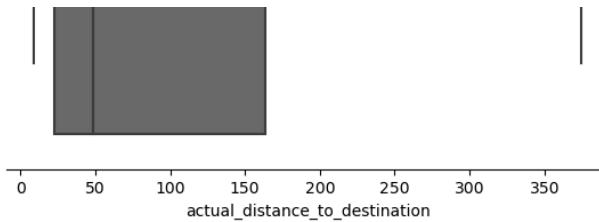
Filtered data of actual_distance_to_destination

actual_distance_to_destination

| | |
|--------------|-------------|
| 0 | 824.732849 |
| 1 | 73.186905 |
| 2 | 1927.404297 |
| 3 | 17.175274 |
| 4 | 127.448502 |
| ... | ... |
| 14782 | 57.762333 |
| 14783 | 15.513784 |
| 14784 | 38.684837 |
| 14785 | 134.723831 |
| 14786 | 66.081528 |

14787 rows × 1 columns

Boxplot of clipped actual_distance_to_destination**Boxplot of filtered actual_distance_to_destination**

**actual_time**

| | |
|--------------|--------|
| 0 | 1562.0 |
| 1 | 143.0 |
| 2 | 3347.0 |
| 3 | 59.0 |
| 4 | 341.0 |
| ... | ... |
| 14782 | 83.0 |
| 14783 | 21.0 |
| 14784 | 282.0 |
| 14785 | 264.0 |
| 14786 | 275.0 |

14787 rows × 1 columns

Clipped data of **actual_time****actual_time**

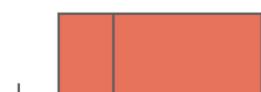
| | |
|--------------|-------|
| 0 | 817.0 |
| 1 | 143.0 |
| 2 | 817.0 |
| 3 | 59.0 |
| 4 | 341.0 |
| ... | ... |
| 14782 | 83.0 |
| 14783 | 21.0 |
| 14784 | 282.0 |
| 14785 | 264.0 |
| 14786 | 275.0 |

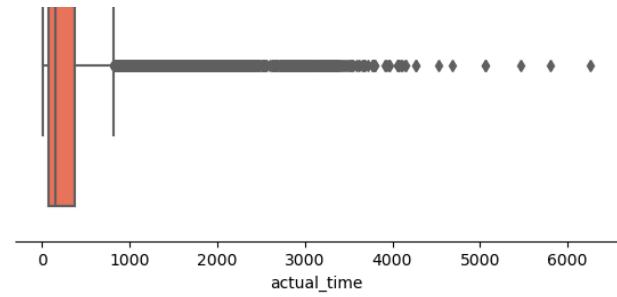
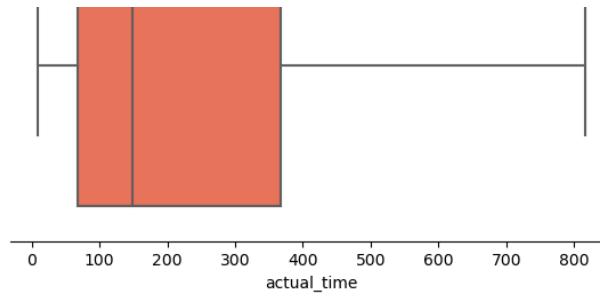
14787 rows × 1 columns

Filtered data of **actual_time****actual_time**

| | |
|--------------|--------|
| 0 | 1562.0 |
| 1 | 143.0 |
| 2 | 3347.0 |
| 3 | 59.0 |
| 4 | 341.0 |
| ... | ... |
| 14782 | 83.0 |
| 14783 | 21.0 |
| 14784 | 282.0 |
| 14785 | 264.0 |
| 14786 | 275.0 |

14787 rows × 1 columns

Boxplot of clipped actual_time**Boxplot of filtered actual_time**

**osrm_time**

| | |
|--------------|--------|
| 0 | 717.0 |
| 1 | 68.0 |
| 2 | 1740.0 |
| 3 | 15.0 |
| 4 | 117.0 |
| ... | ... |
| 14782 | 62.0 |
| 14783 | 12.0 |
| 14784 | 48.0 |
| 14785 | 179.0 |
| 14786 | 68.0 |

14787 rows × 1 columns

Clipped data of osrm_time

osrm_time

| | |
|--------------|-------|
| 0 | 376.5 |
| 1 | 68.0 |
| 2 | 376.5 |
| 3 | 15.0 |
| 4 | 117.0 |
| ... | ... |
| 14782 | 62.0 |
| 14783 | 12.0 |
| 14784 | 48.0 |
| 14785 | 179.0 |
| 14786 | 68.0 |

14787 rows × 1 columns

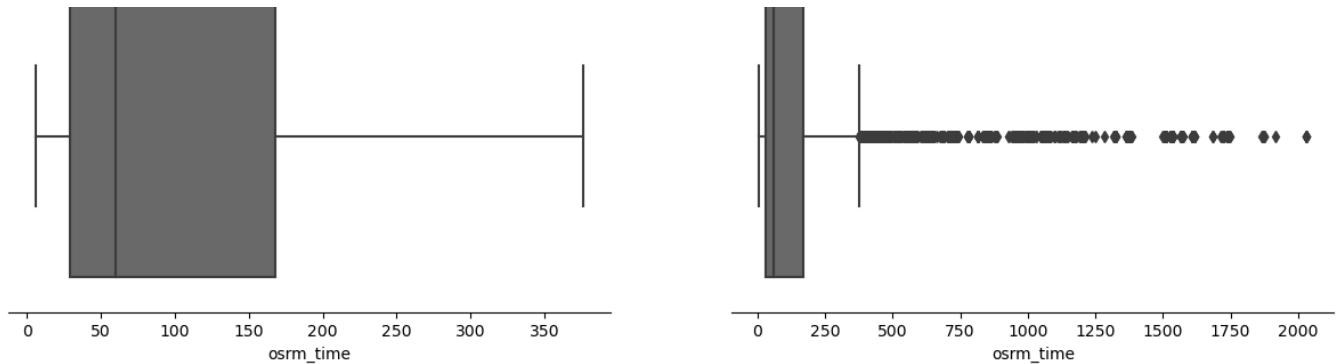
Filtered data of osrm_time

osrm_time

| | |
|--------------|--------|
| 0 | 717.0 |
| 1 | 68.0 |
| 2 | 1740.0 |
| 3 | 15.0 |
| 4 | 117.0 |
| ... | ... |
| 14782 | 62.0 |
| 14783 | 12.0 |
| 14784 | 48.0 |
| 14785 | 179.0 |
| 14786 | 68.0 |

14787 rows × 1 columns

Boxplot of clipped osrm_time**Boxplot of filtered osrm_time**

**osrm_distance**

| | |
|--------------|-------------|
| 0 | 991.352295 |
| 1 | 85.111000 |
| 2 | 2354.066650 |
| 3 | 19.680000 |
| 4 | 146.791794 |
| ... | ... |
| 14782 | 73.462997 |
| 14783 | 16.088200 |
| 14784 | 58.903702 |
| 14785 | 171.110306 |
| 14786 | 80.578705 |

14787 rows × 1 columns

Clipped data of osrm_distance

osrm_distance

| | |
|--------------|------------|
| 0 | 470.475158 |
| 1 | 85.111000 |
| 2 | 470.475158 |
| 3 | 19.680000 |
| 4 | 146.791794 |
| ... | ... |
| 14782 | 73.462997 |
| 14783 | 16.088200 |
| 14784 | 58.903702 |
| 14785 | 171.110306 |
| 14786 | 80.578705 |

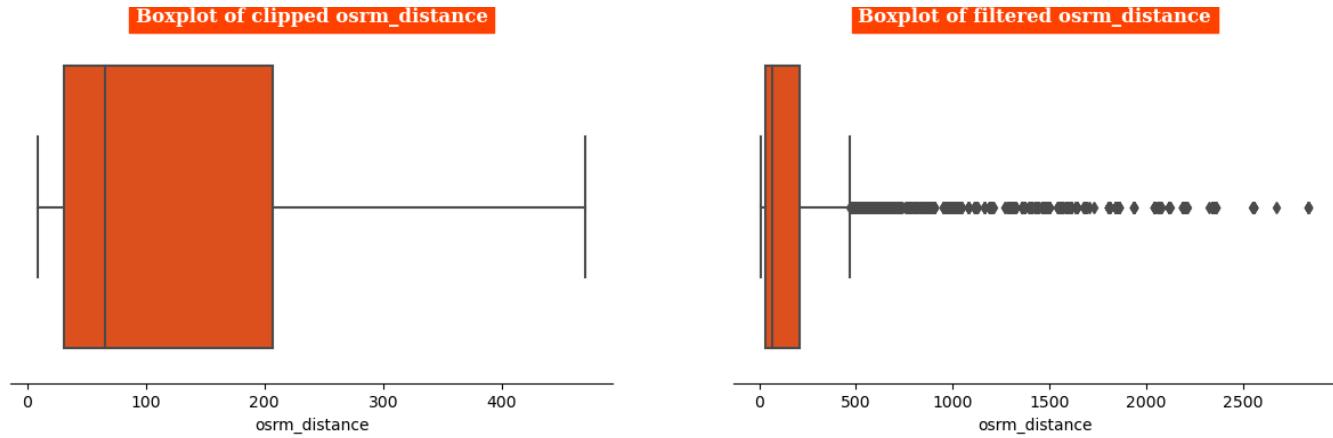
14787 rows × 1 columns

Filtered data of osrm_distance

osrm_distance

| | |
|--------------|-------------|
| 0 | 991.352295 |
| 1 | 85.111000 |
| 2 | 2354.066650 |
| 3 | 19.680000 |
| 4 | 146.791794 |
| ... | ... |
| 14782 | 73.462997 |
| 14783 | 16.088200 |
| 14784 | 58.903702 |
| 14785 | 171.110306 |
| 14786 | 80.578705 |

14787 rows × 1 columns



segment_actual_time

| | |
|--------------|--------|
| 0 | 1548.0 |
| 1 | 141.0 |
| 2 | 3308.0 |
| 3 | 59.0 |
| 4 | 340.0 |
| ... | ... |
| 14782 | 82.0 |
| 14783 | 21.0 |
| 14784 | 281.0 |
| 14785 | 258.0 |
| 14786 | 274.0 |

14787 rows × 1 columns

Clipped data of segment_actual_time

segment_actual_time

| | |
|--------------|-------|
| 0 | 811.0 |
| 1 | 141.0 |
| 2 | 811.0 |
| 3 | 59.0 |
| 4 | 340.0 |
| ... | ... |
| 14782 | 82.0 |
| 14783 | 21.0 |
| 14784 | 281.0 |
| 14785 | 258.0 |
| 14786 | 274.0 |

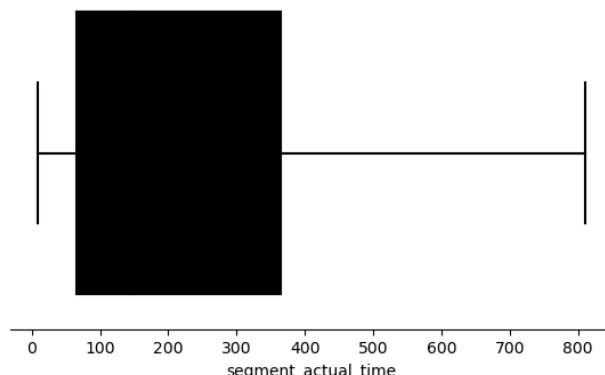
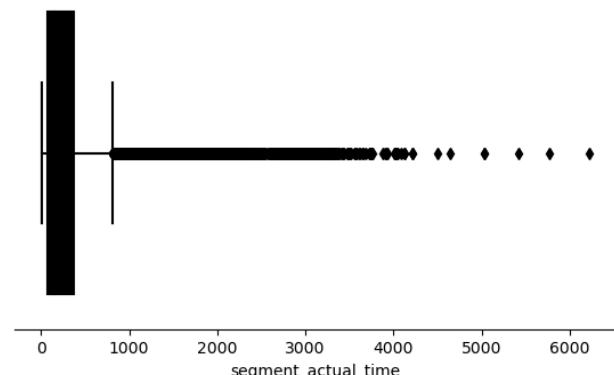
14787 rows × 1 columns

Filtered data of segment_actual_time

segment_actual_time

| | |
|--------------|--------|
| 0 | 1548.0 |
| 1 | 141.0 |
| 2 | 3308.0 |
| 3 | 59.0 |
| 4 | 340.0 |
| ... | ... |
| 14782 | 82.0 |
| 14783 | 21.0 |
| 14784 | 281.0 |
| 14785 | 258.0 |
| 14786 | 274.0 |

14787 rows × 1 columns

Boxplot of clipped segment_actual_time**Boxplot of filtered segment_actual_time****segment_osrm_time**

| | |
|--------------|--------|
| 0 | 1008.0 |
| 1 | 65.0 |
| 2 | 1941.0 |
| 3 | 16.0 |
| 4 | 115.0 |
| ... | ... |
| 14782 | 62.0 |
| 14783 | 11.0 |
| 14784 | 88.0 |
| 14785 | 221.0 |
| 14786 | 67.0 |

14787 rows × 1 columns

Clipped data of segment_osrm_time

segment_osrm_time

| | |
|--------------|-------|
| 0 | 415.0 |
| 1 | 65.0 |
| 2 | 415.0 |
| 3 | 16.0 |
| 4 | 115.0 |
| ... | ... |
| 14782 | 62.0 |
| 14783 | 11.0 |
| 14784 | 88.0 |
| 14785 | 221.0 |
| 14786 | 67.0 |

14787 rows × 1 columns

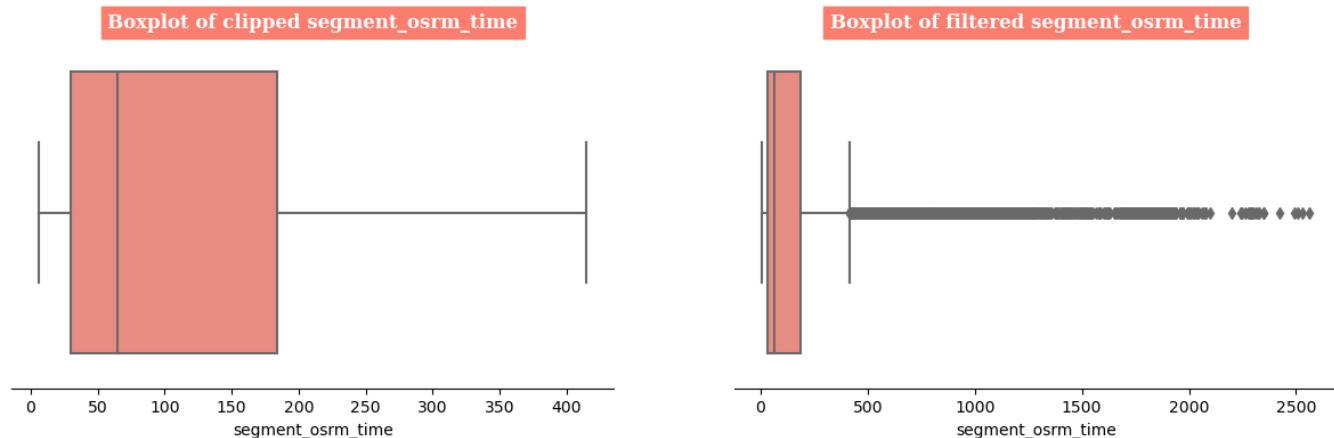
Filtered data of segment_osrm_time

segment_osrm_time

| | |
|--------------|--------|
| 0 | 1008.0 |
| 1 | 65.0 |
| 2 | 1941.0 |
| 3 | 16.0 |
| 4 | 115.0 |
| ... | ... |
| 14782 | 62.0 |
| 14783 | 11.0 |
| 14784 | 88.0 |

| | |
|--------------|-------|
| 14785 | 221.0 |
| 14786 | 67.0 |

14787 rows × 1 columns

**segment_osrm_distance**

| | |
|--------------|-------------|
| 0 | 1320.473267 |
| 1 | 84.189400 |
| 2 | 2545.267822 |
| 3 | 19.876600 |
| 4 | 146.791901 |
| ... | ... |
| 14782 | 64.855103 |
| 14783 | 16.088299 |
| 14784 | 104.886597 |
| 14785 | 223.532394 |
| 14786 | 80.578705 |

14787 rows × 1 columns

Clipped data of segment_osrm_distance

segment_osrm_distance

| | |
|--------------|------------|
| 0 | 492.533245 |
| 1 | 84.189400 |
| 2 | 492.533245 |
| 3 | 19.876600 |
| 4 | 146.791901 |
| ... | ... |
| 14782 | 64.855103 |
| 14783 | 16.088299 |
| 14784 | 104.886597 |
| 14785 | 223.532394 |
| 14786 | 80.578705 |

14787 rows × 1 columns

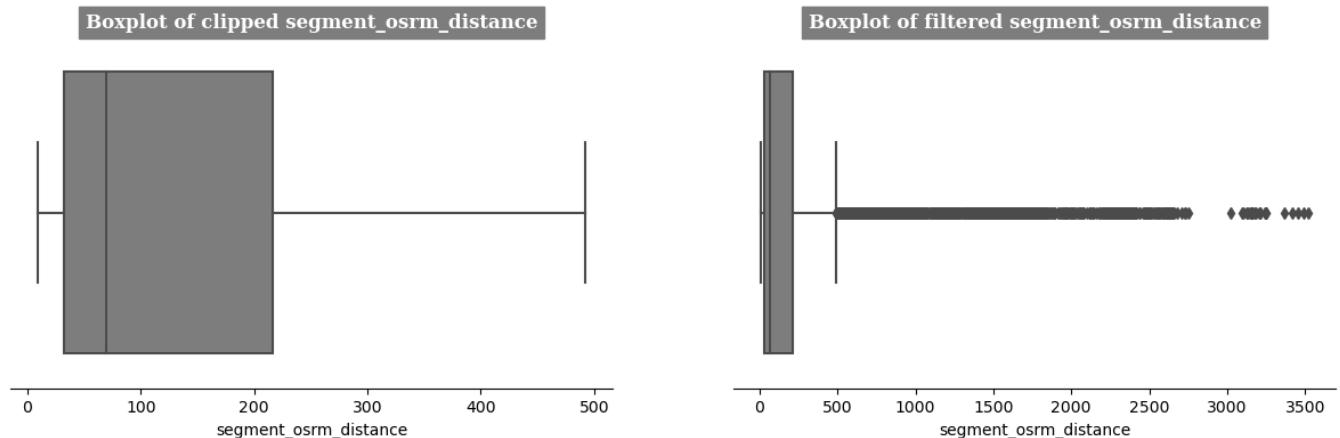
Filtered data of segment_osrm_distance

segment_osrm_distance

| | |
|--------------|-------------|
| 0 | 1320.473267 |
| 1 | 84.189400 |
| 2 | 2545.267822 |
| 3 | 19.876600 |
| 4 | 146.791901 |
| ... | ... |
| 14782 | 64.855103 |
| 14783 | 16.088299 |

| | |
|--------------|------------|
| 14784 | 104.886597 |
| 14785 | 223.532394 |
| 14786 | 80.578705 |

14787 rows × 1 columns

**segment_actual_time_sum**

| | |
|--------------|--------|
| 0 | 1548.0 |
| 1 | 141.0 |
| 2 | 3308.0 |
| 3 | 59.0 |
| 4 | 340.0 |
| ... | ... |
| 14782 | 82.0 |
| 14783 | 21.0 |
| 14784 | 281.0 |
| 14785 | 258.0 |
| 14786 | 274.0 |

14787 rows × 1 columns

Clipped data of segment_actual_time_sum

| | segment_actual_time_sum |
|--------------|-------------------------|
| 0 | 811.0 |
| 1 | 141.0 |
| 2 | 811.0 |
| 3 | 59.0 |
| 4 | 340.0 |
| ... | ... |
| 14782 | 82.0 |
| 14783 | 21.0 |
| 14784 | 281.0 |
| 14785 | 258.0 |
| 14786 | 274.0 |

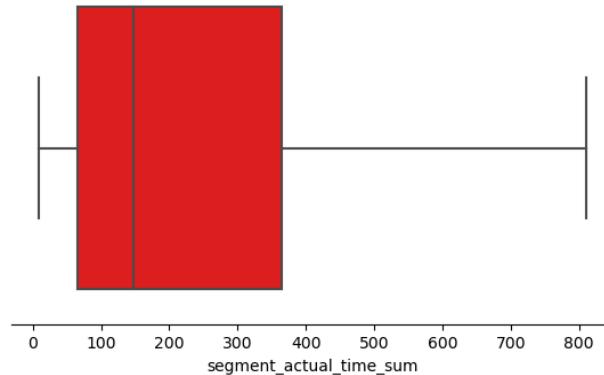
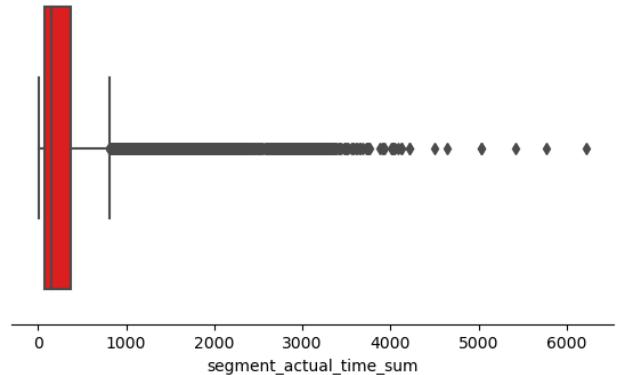
14787 rows × 1 columns

Filtered data of segment_actual_time_sum

| | segment_actual_time_sum |
|----------|-------------------------|
| 0 | 1548.0 |
| 1 | 141.0 |
| 2 | 3308.0 |
| 3 | 59.0 |
| 4 | 340.0 |
| ... | ... |

| | |
|-------|-------|
| 14782 | 82.0 |
| 14783 | 21.0 |
| 14784 | 281.0 |
| 14785 | 258.0 |
| 14786 | 274.0 |

14787 rows × 1 columns

Boxplot of clipped segment_actual_time_sum**Boxplot of filtered segment_actual_time_sum****segment_osrm_time_sum**

| | |
|-------|--------|
| 0 | 1008.0 |
| 1 | 65.0 |
| 2 | 1941.0 |
| 3 | 16.0 |
| 4 | 115.0 |
| ... | ... |
| 14782 | 62.0 |
| 14783 | 11.0 |
| 14784 | 88.0 |
| 14785 | 221.0 |
| 14786 | 67.0 |

14787 rows × 1 columns

Clipped data of segment_osrm_time_sum

segment_osrm_time_sum

| | |
|-------|-------|
| 0 | 415.0 |
| 1 | 65.0 |
| 2 | 415.0 |
| 3 | 16.0 |
| 4 | 115.0 |
| ... | ... |
| 14782 | 62.0 |
| 14783 | 11.0 |
| 14784 | 88.0 |
| 14785 | 221.0 |
| 14786 | 67.0 |

14787 rows × 1 columns

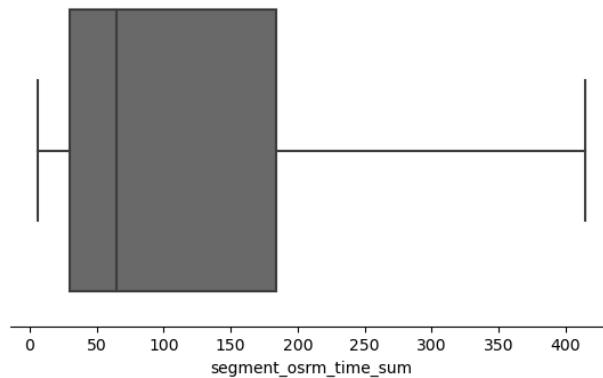
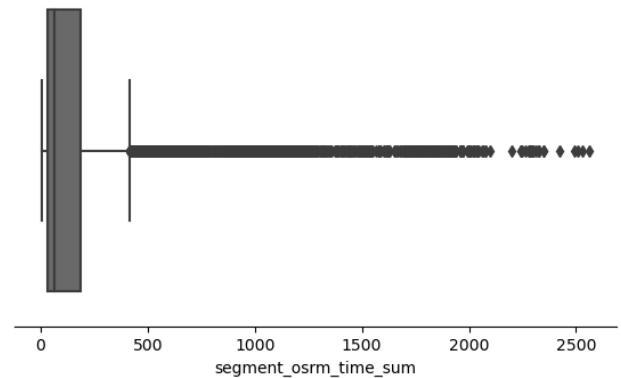
Filtered data of segment_osrm_time_sum

segment_osrm_time_sum

| | |
|---|--------|
| 0 | 1008.0 |
| 1 | 65.0 |
| 2 | 1941.0 |
| 3 | 16.0 |

| | |
|--------------|-------|
| 4 | 115.0 |
| ... | ... |
| 14782 | 62.0 |
| 14783 | 11.0 |
| 14784 | 88.0 |
| 14785 | 221.0 |
| 14786 | 67.0 |

14787 rows × 1 columns

Boxplot of clipped segment_osrm_time_sum**Boxplot of filtered segment_osrm_time_sum****segment_osrm_distance_sum**

| | |
|--------------|-------------|
| 0 | 1320.473267 |
| 1 | 84.189400 |
| 2 | 2545.267822 |
| 3 | 19.876600 |
| 4 | 146.791901 |
| ... | ... |
| 14782 | 64.855103 |
| 14783 | 16.088299 |
| 14784 | 104.886597 |
| 14785 | 223.532394 |
| 14786 | 80.578705 |

14787 rows × 1 columns

Clipped data of segment_osrm_distance_sum

segment_osrm_distance_sum

| | |
|--------------|------------|
| 0 | 492.533245 |
| 1 | 84.189400 |
| 2 | 492.533245 |
| 3 | 19.876600 |
| 4 | 146.791901 |
| ... | ... |
| 14782 | 64.855103 |
| 14783 | 16.088299 |
| 14784 | 104.886597 |
| 14785 | 223.532394 |
| 14786 | 80.578705 |

14787 rows × 1 columns

Filtered data of segment_osrm_distance_sum

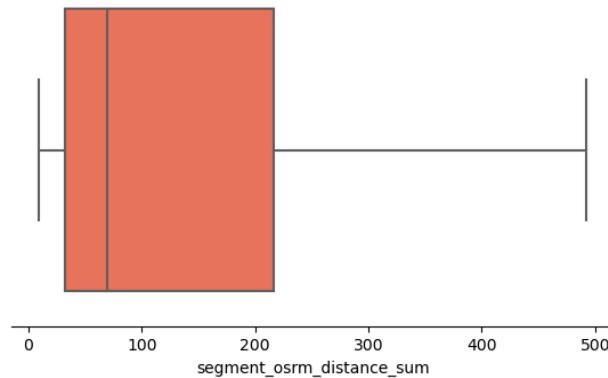
segment_osrm_distance_sum

| | |
|---|-------------|
| 0 | 1320.473267 |
| 1 | 84.189400 |
| 2 | 2545.267822 |

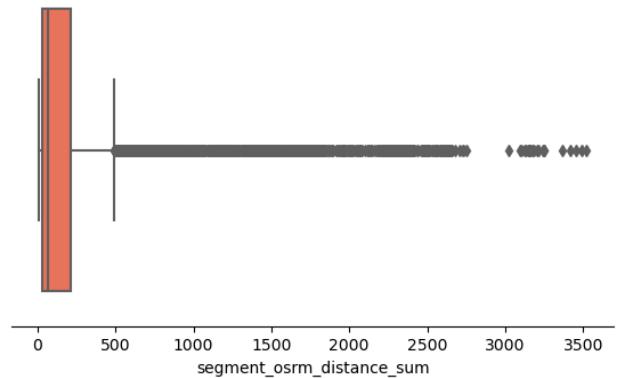
```
2043.201022  
3      19.876600  
4     146.791901  
...  
14782    64.855103  
14783   16.088299  
14784  104.886597  
14785 223.532394  
14786  80.578705
```

14787 rows × 1 columns

Boxplot of clipped segment_osrm_distance_sum



Boxplot of filtered segment_osrm_distance_sum



⌄ ⚡ Understanding:

- Here we see that the data after removing outliers has outliers. It has to be understood that q1 and q3 dont have to be always 25th percentile and 75th percentile. **Try changing q1 and q3 to 10th percentile to 90th percentile and plot and see...**
- Clipped data replaces the outlier values with specified values.
- Here, I have proceeded with both clipped and filtered data(with reduced outliers) for further analysis.

```
1 num_df = numerical_columns.copy()
2 num_df
```

| | od_time_diff_hour | start_scan_to_end_scan | actual_distance_to_destination | actual_time | osrm_time | osrm_distance | segment_act |
|-------|-------------------|------------------------|--------------------------------|-------------|-----------|---------------|-------------|
| 0 | 37.668497 | 2259.0 | 824.732849 | 1562.0 | 717.0 | 991.352295 | |
| 1 | 3.026865 | 180.0 | 73.186905 | 143.0 | 68.0 | 85.111000 | |
| 2 | 65.572709 | 3933.0 | 1927.404297 | 3347.0 | 1740.0 | 2354.066650 | |
| 3 | 1.674916 | 100.0 | 17.175274 | 59.0 | 15.0 | 19.680000 | |
| 4 | 11.972484 | 717.0 | 127.448502 | 341.0 | 117.0 | 146.791794 | |
| ... | ... | ... | ... | ... | ... | ... | ... |
| 14782 | 4.300482 | 257.0 | 57.762333 | 83.0 | 62.0 | 73.462997 | |
| 14783 | 1.009842 | 60.0 | 15.513784 | 21.0 | 12.0 | 16.088200 | |
| 14784 | 7.035331 | 421.0 | 38.684837 | 282.0 | 48.0 | 58.903702 | |
| 14785 | 5.808548 | 347.0 | 134.723831 | 264.0 | 179.0 | 171.110306 | |
| 14786 | 5.906793 | 353.0 | 66.081528 | 275.0 | 68.0 | 80.578705 | |

14787 rows × 12 columns

```
1 num_cols
```

```
2 ⌄ ['od_time_diff_hour',
3   'start_scan_to_end_scan',
4   'actual_distance_to_destination',
5   'actual_time',
6   'osrm_time',
7   'osrm_distance',
8   'segment_actual_time',
9   'segment_osrm_time',
10  'segment_osrm_distance',
11  'segment_actual_time_sum',
12  'segment_osrm_time_sum',
13  'segment_osrm_distance_sum']
```

```
1 Q1 = np.percentile(num_df[num_cols], 25)
2 Q3 = np.percentile(num_df[num_cols], 75)
3 IQR = Q3 - Q1
4
5 lower_bound = Q1 - (1.5 * IQR)
6 upper_bound = Q3 + (1.5 * IQR)
7
8 clipped_num_df = np.clip(num_df[num_cols], lower_bound, upper_bound)
9 display(clipped_num_df)
10
11 filtered_num_df = num_df[num_cols][(num_df[num_cols] >= lower_bound) | (num_df[num_cols] <= upper_bound)]
12 display(filtered_num_df)
```

| | od_time_diff_hour | start_scan_to_end_scan | actual_distance_to_destination | actual_time | osrm_time | osrm_distance | segment_ac |
|-------|-------------------|------------------------|--------------------------------|-------------|------------|---------------|------------|
| 0 | 37.668497 | 543.285302 | | 543.285302 | 543.285302 | 543.285302 | 5 |
| 1 | 3.026865 | 180.000000 | | 73.186905 | 143.000000 | 68.000000 | 85.111000 |
| 2 | 65.572709 | 543.285302 | | 543.285302 | 543.285302 | 543.285302 | 5 |
| 3 | 1.674916 | 100.000000 | | 17.175274 | 59.000000 | 15.000000 | 19.680000 |
| 4 | 11.972484 | 543.285302 | | 127.448502 | 341.000000 | 117.000000 | 146.791794 |
| ... | ... | ... | | ... | ... | ... | ... |
| 14782 | 4.300482 | 257.000000 | | 57.762333 | 83.000000 | 62.000000 | 73.462997 |
| 14783 | 1.009842 | 60.000000 | | 15.513784 | 21.000000 | 12.000000 | 16.088200 |
| 14784 | 7.035331 | 421.000000 | | 38.684837 | 282.000000 | 48.000000 | 58.903702 |
| 14785 | 5.808548 | 347.000000 | | 134.723831 | 264.000000 | 179.000000 | 171.110306 |
| 14786 | 5.906793 | 353.000000 | | 66.081528 | 275.000000 | 68.000000 | 80.578705 |

14787 rows × 12 columns

| | od_time_diff_hour | start_scan_to_end_scan | actual_distance_to_destination | actual_time | osrm_time | osrm_distance | segment_act |
|-------|-------------------|------------------------|--------------------------------|-------------|-----------|---------------|-------------|
| 0 | 37.668497 | 2259.0 | | 824.732849 | 1562.0 | 717.0 | 991.352295 |
| 1 | 3.026865 | 180.0 | | 73.186905 | 143.0 | 68.0 | 85.111000 |
| 2 | 65.572709 | 3933.0 | | 1927.404297 | 3347.0 | 1740.0 | 2354.066650 |
| 3 | 1.674916 | 100.0 | | 17.175274 | 59.0 | 15.0 | 19.680000 |
| 4 | 11.972484 | 717.0 | | 127.448502 | 341.0 | 117.0 | 146.791794 |
| ... | ... | ... | | ... | ... | ... | ... |
| 14782 | 4.300482 | 257.0 | | 57.762333 | 83.0 | 62.0 | 73.462997 |
| 14783 | 1.009842 | 60.0 | | 15.513784 | 21.0 | 12.0 | 16.088200 |
| 14784 | 7.035331 | 421.0 | | 38.684837 | 282.0 | 48.0 | 58.903702 |
| 14785 | 5.808548 | 347.0 | | 134.723831 | 264.0 | 179.0 | 171.110306 |
| 14786 | 5.906793 | 353.0 | | 66.081528 | 275.0 | 68.0 | 80.578705 |

14787 rows × 12 columns

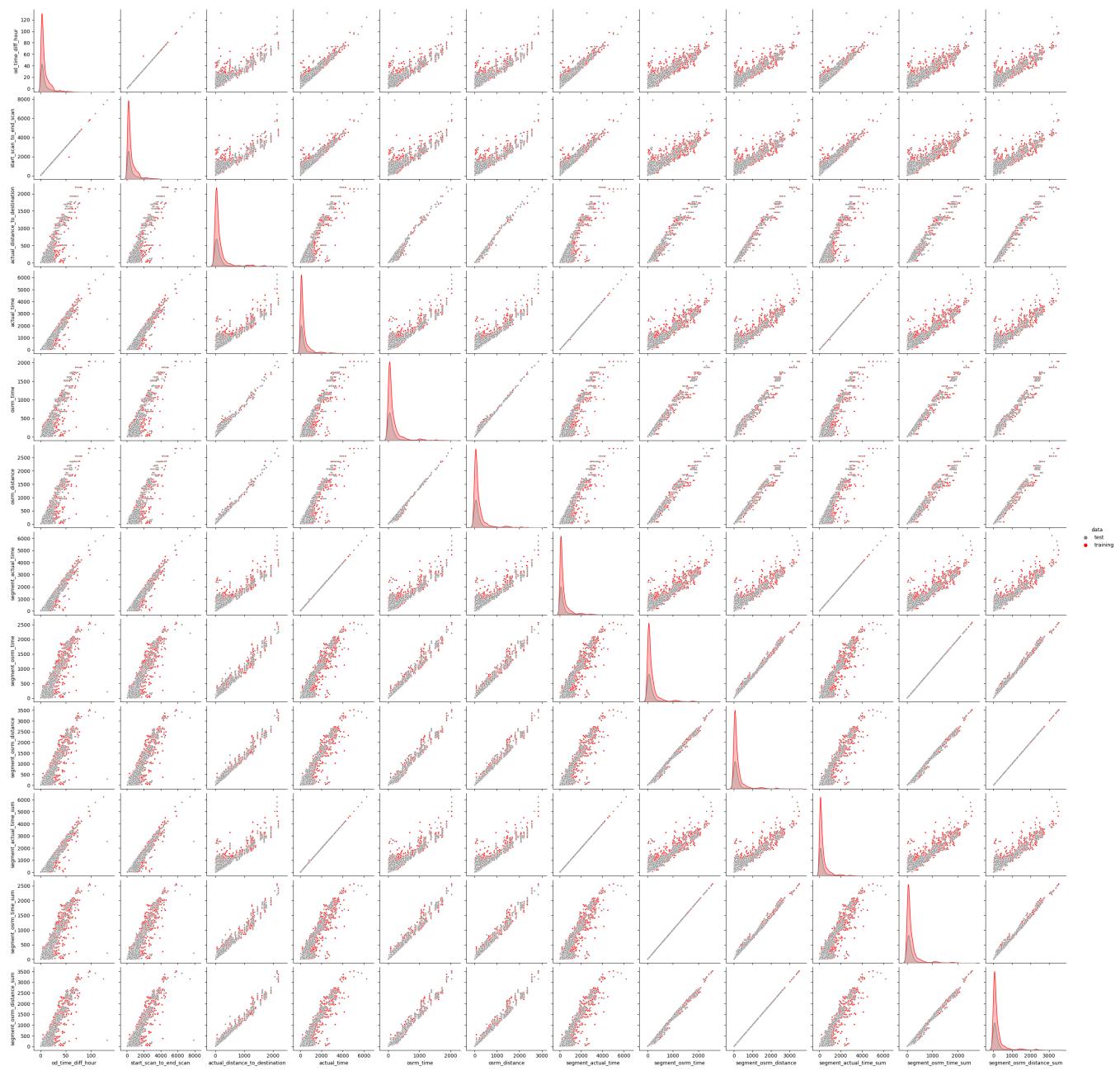
```

1 plt.figure(figsize=(14,0.05))
2 plt.axis('off')
3 plt.title(f' Pairplot Analysis',fontfamily='serif',fontweight='bold',fontsize=15,backgroundcolor='crimson',color='w')
4 sns.pairplot(data = trip_df,vars = num_cols,hue='data',markers = '.',palette=cp)
5 plt.show()

```



Pairplot Analysis



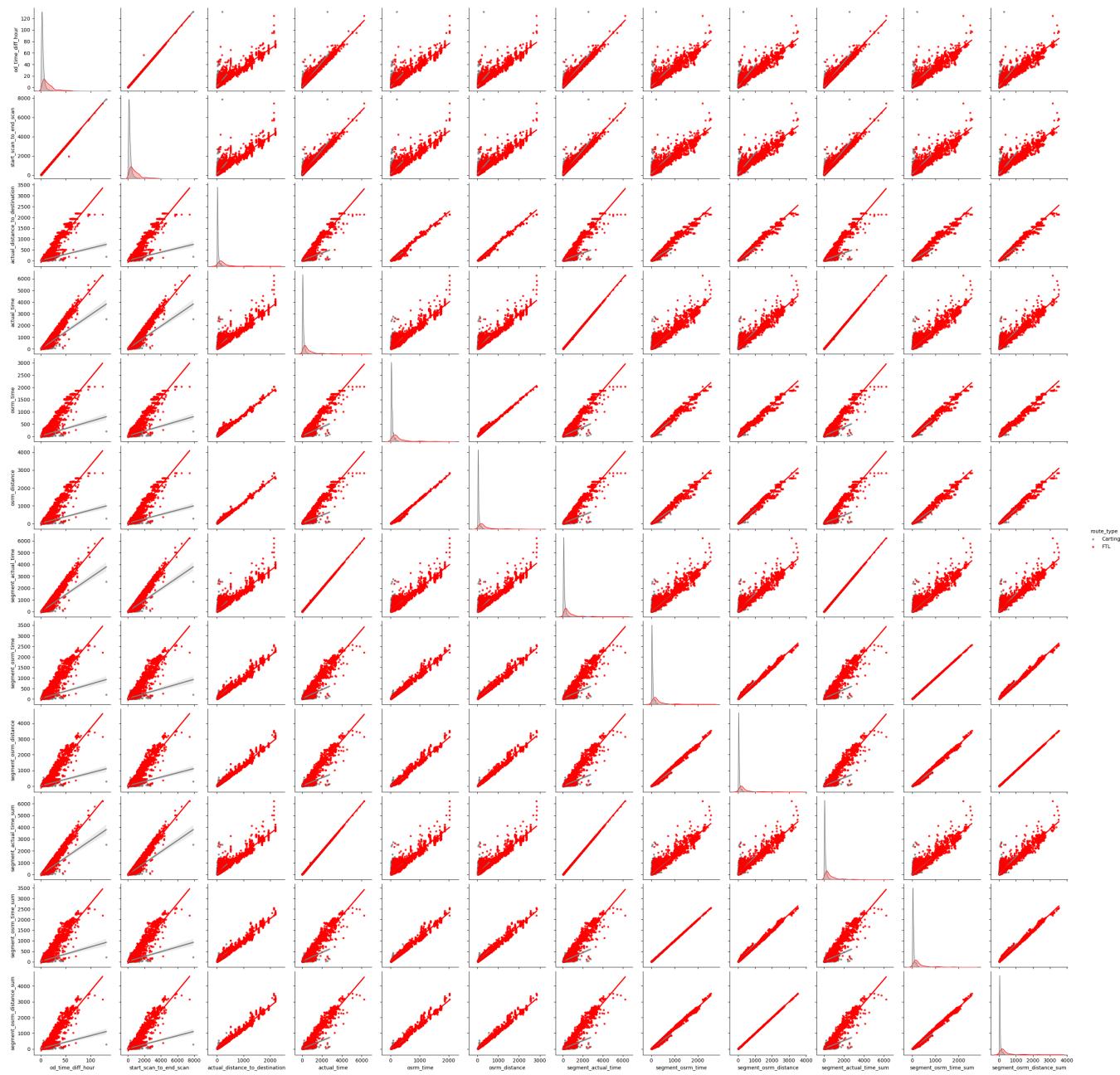
```

1 plt.figure(figsize=(14,0.05))
2 plt.axis('off')
3 plt.title(f' Pairplot Analysis',fontfamily='serif',fontweight='bold',fontsize=15,backgroundcolor='dimgrey',color='w')
4 sns.pairplot(data = trip_df,vars=num_cols,kind = 'reg',hue='route_type',markers = '.',palette=cp)
5 plt.show()

```



Pairplot Analysis



```

1 clipped_df_corr = clipped_num_df.corr()
2 clipped_df_corr
3
4 filtered_df_corr = filtered_num_df.corr()
5 filtered_df_corr

```