

```
In [ ]: #importing essential libraries
```

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
```

```
In [ ]: #Load CSV file into dataframe
```

```
df=pd.read_csv('delhivery_data.csv')
```

```
In [ ]: #having Look of dataset for first 5 rows
```

```
df.head()
```

```
Out[ ]:
```

	data	trip_creation_time	route_schedule_uuid	route_type	trip_uuid	source_cer
0	training	2018-09-20 02:35:36.476840	thanos::sroute:eb7bfc78- b351-4c0e-a951- fa3d5c3...	Carting	153741093647649320	trip- IND388121/
1	training	2018-09-20 02:35:36.476840	thanos::sroute:eb7bfc78- b351-4c0e-a951- fa3d5c3...	Carting	153741093647649320	trip- IND388121/
2	training	2018-09-20 02:35:36.476840	thanos::sroute:eb7bfc78- b351-4c0e-a951- fa3d5c3...	Carting	153741093647649320	trip- IND388121/
3	training	2018-09-20 02:35:36.476840	thanos::sroute:eb7bfc78- b351-4c0e-a951- fa3d5c3...	Carting	153741093647649320	trip- IND388121/
4	training	2018-09-20 02:35:36.476840	thanos::sroute:eb7bfc78- b351-4c0e-a951- fa3d5c3...	Carting	153741093647649320	trip- IND388121/

5 rows × 24 columns

```
In [ ]: #gathering more information about data
```

```
df.info()
```

```

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

```

In [ ]: `#checking shape of dataset`

```
df.shape
```

Out[ ]: (144867, 24)

In [ ]: `#statistical info about dataset having numeric columns`

```
df.describe()
```

Out[ ]:

	<b>start_scan_to_end_scan</b>	<b>cutoff_factor</b>	<b>actual_distance_to_destination</b>	<b>actual_time</b>	<b>osr</b>
<b>count</b>	144867.000000	144867.000000	144867.000000	144867.000000	144867
<b>mean</b>	961.262986	232.926567	234.073372	416.927527	213
<b>std</b>	1037.012769	344.755577	344.990009	598.103621	308
<b>min</b>	20.000000	9.000000	9.000045	9.000000	6
<b>25%</b>	161.000000	22.000000	23.355874	51.000000	27
<b>50%</b>	449.000000	66.000000	66.126571	132.000000	64
<b>75%</b>	1634.000000	286.000000	286.708875	513.000000	257
<b>max</b>	7898.000000	1927.000000	1927.447705	4532.000000	1686

In [ ]: `#statistical info about dataset having objective coulmns and transposing it`

```
df.describe(include='object').T
```

Out[ ]:

		count	unique	top	freq
	<b>data</b>	144867	2		training 104858
	<b>trip_creation_time</b>	144867	14817	2018-09-28 05:23:15.359220	101
	<b>route_schedule_uuid</b>	144867	1504	thanos::sroute:4029a8a2-6c74-4b7e-a6d8-f9e069f...	1812
	<b>route_type</b>	144867	2		FTL 99660
	<b>trip_uuid</b>	144867	14817	trip-153811219535896559	101
	<b>source_center</b>	144867	1508	IND00000ACB	23347
	<b>source_name</b>	144574	1498	Gurgaon_Bilaspur_HB (Haryana)	23347
	<b>destination_center</b>	144867	1481	IND00000ACB	15192
	<b>destination_name</b>	144606	1468	Gurgaon_Bilaspur_HB (Haryana)	15192
	<b>od_start_time</b>	144867	26369	2018-09-21 18:37:09.322207	81
	<b>od_end_time</b>	144867	26369	2018-09-24 09:59:15.691618	81
	<b>cutoff_timestamp</b>	144867	93180	2018-09-24 05:19:20	40

In [ ]: #checking for null values in percentage

df.isnull().sum()/len(df)\*100

Out[ ]:

```

data                      0.000000
trip_creation_time        0.000000
route_schedule_uuid        0.000000
route_type                 0.000000
trip_uuid                  0.000000
source_center               0.000000
source_name                 0.202254
destination_center          0.000000
destination_name            0.180165
od_start_time                0.000000
od_end_time                  0.000000
start_scan_to_end_scan      0.000000
is_cutoff                   0.000000
cutoff_factor                0.000000
cutoff_timestamp              0.000000
actual_distance_to_destination 0.000000
actual_time                  0.000000
osrm_time                    0.000000
osrm_distance                 0.000000
factor                       0.000000
segment_actual_time          0.000000
segment_osrm_time             0.000000
segment_osrm_distance         0.000000
segment_factor                 0.000000
dtype: float64

```

In [ ]: #Changing data type of feature

```

df["trip_creation_time"] = pd.to_datetime(df["trip_creation_time"])
df["od_start_time"] = pd.to_datetime(df["od_start_time"])
df["od_end_time"] = pd.to_datetime(df["od_end_time"])

```

In [ ]: #Checking range of data

df["trip\_creation\_time"].dt.month\_name().value\_counts()

```
Out[ ]: September      127349
          October       17518
          Name: trip_creation_time, dtype: int64
```

```
In [ ]: df["trip_creation_time"].dt.year.value_counts()
```

```
Out[ ]: 2018      144867
          Name: trip_creation_time, dtype: int64
```

Data are ranging from the month of September-2018 to October-2018.

```
In [ ]: #Checking no. of Unique Categories of Features
df.nunique()
```

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

There are total 14817 different trips of data available.

There are 1508 unique source\_center.

There are 1481 unique destination\_center.

There are total 1504 delivery routes.

```
In [ ]: #Univariate Analysis

num_vars = df.select_dtypes(include=np.number).columns.tolist()

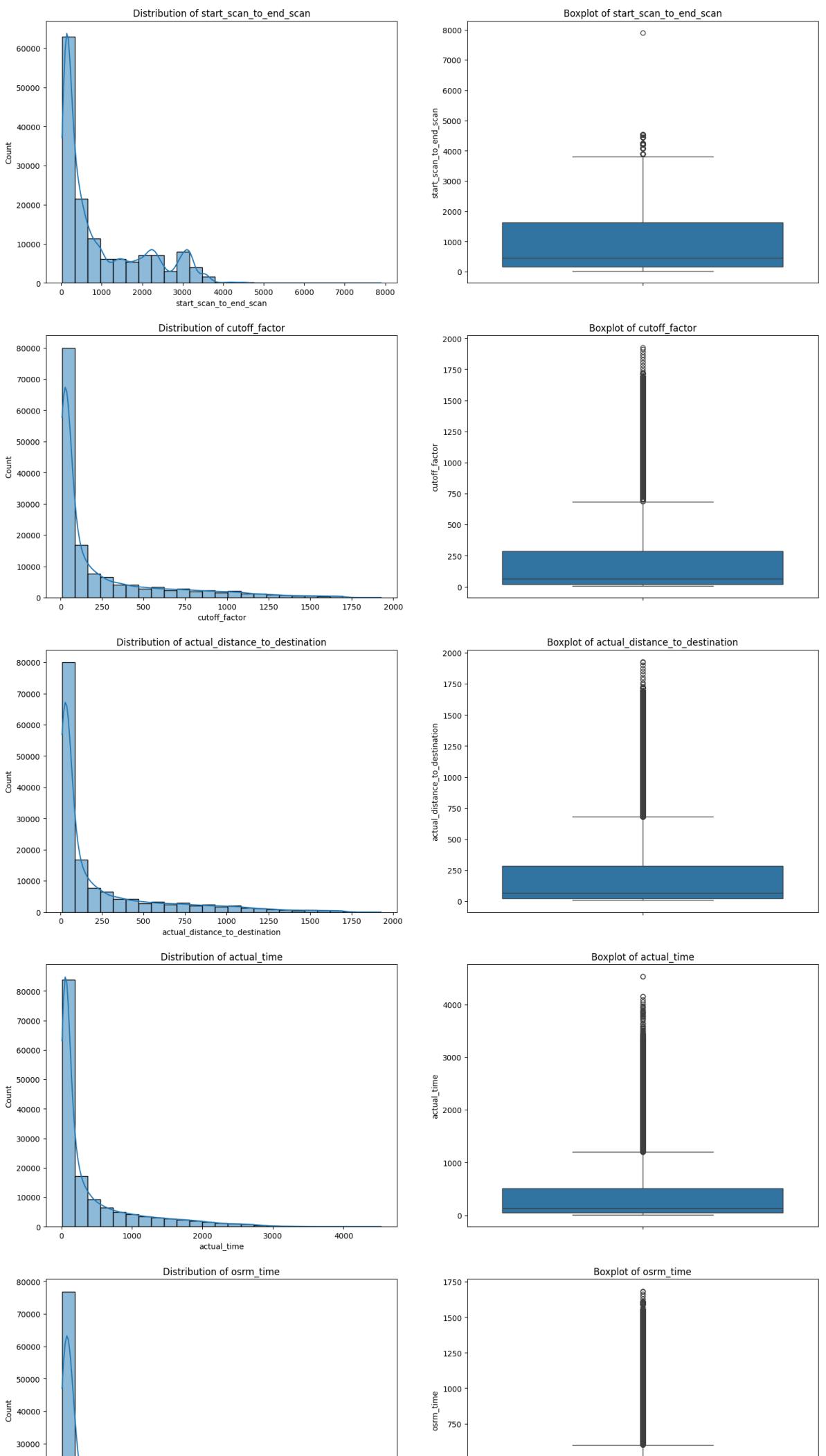
fig, ax = plt.subplots(nrows=11, ncols=2, figsize=(18, 80))

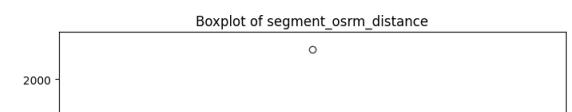
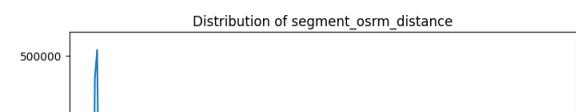
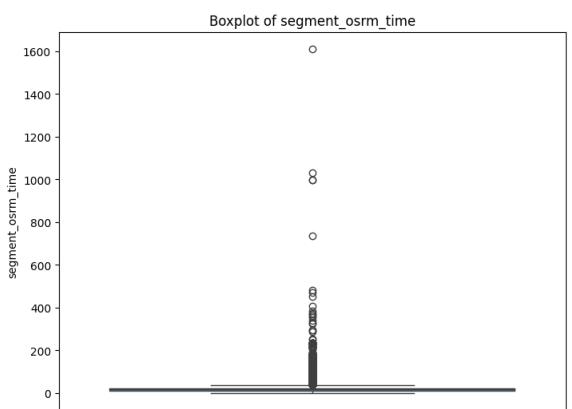
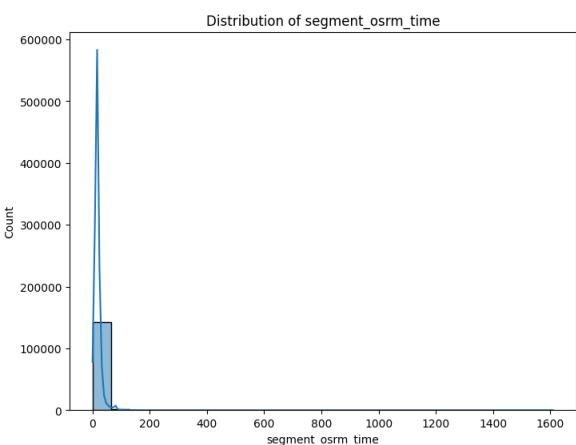
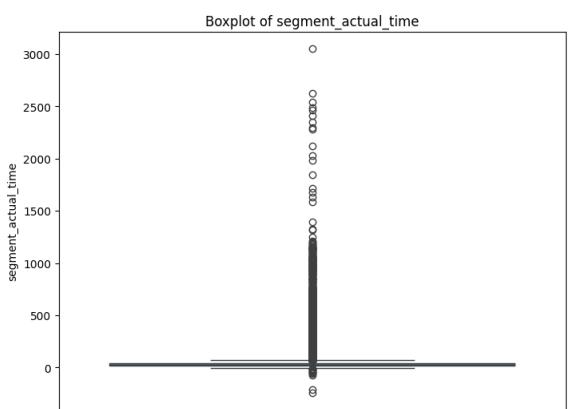
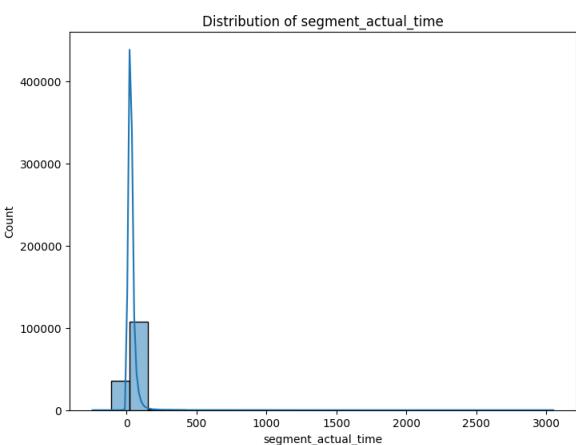
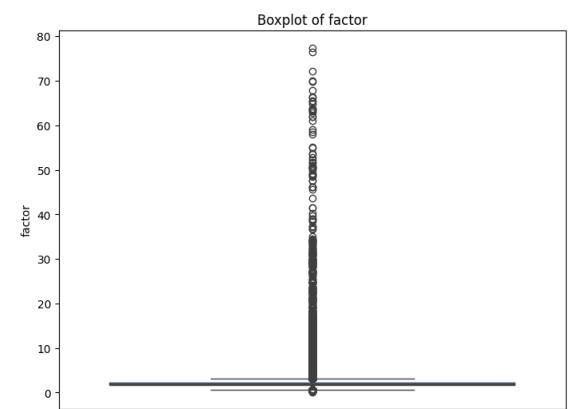
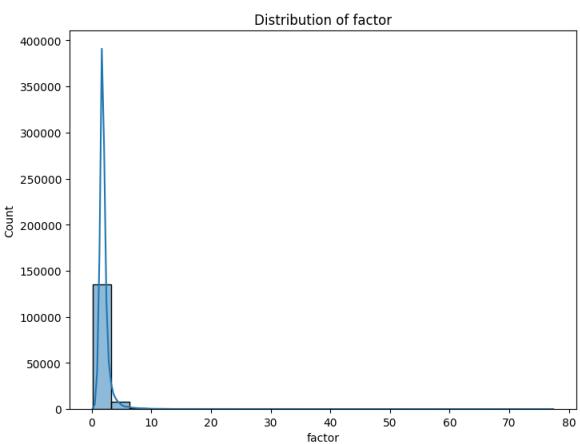
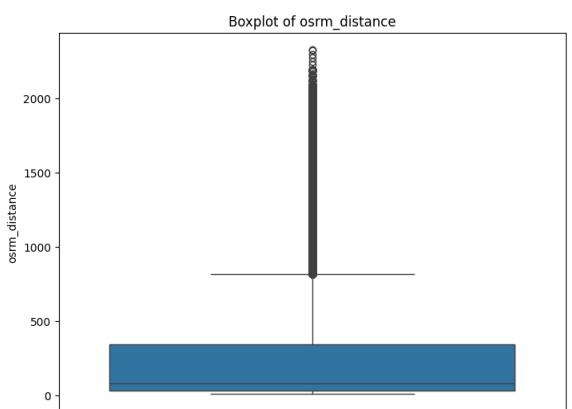
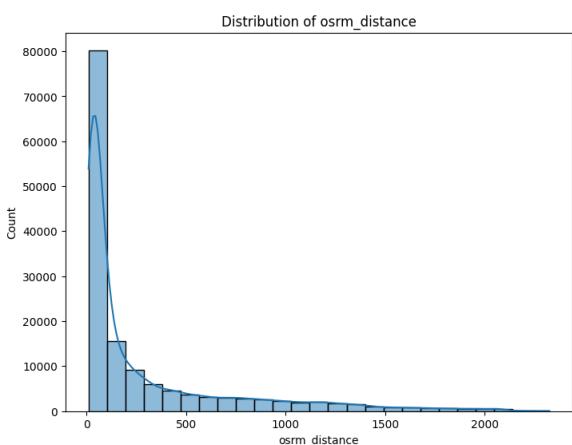
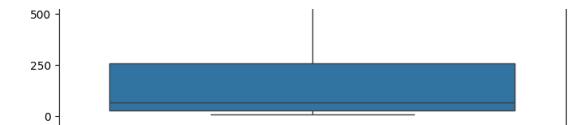
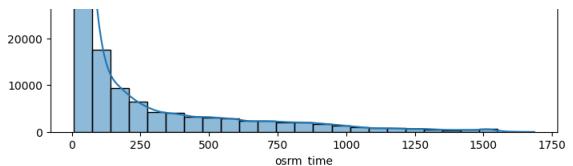
for i in range(len(num_vars)):

    sns.histplot(x=df[num_vars[i]], kde=True, bins = 25, ax=ax[i, 0])
    ax[i, 0].set_title(f"Distribution of {num_vars[i]}")

    sns.boxplot(y = df[num_vars[i]], ax=ax[i, 1], data=df)
    ax[i, 1].set_title(f"Boxplot of {num_vars[i]}")
```

```
plt.show()
```





```
In [ ]: #Extracting Features by splitting columns into multiple columns

df["source_city"] = df["source_name"].str.split(" ",n=1,expand=True)[0].str.split('
df["source_state"] = df["source_name"].str.split(" ",n=1,expand=True)[1].str.replace('

df["destination_city"] = df["destination_name"].str.split(" ",n=1,expand=True)[0].s
df["destination_state"] = df["destination_name"].str.split(" ",n=1,expand=True)[1].s

df["source_pincode"] = df["source_center"].apply(lambda x : x[3:9] )
df["destination_pincode"] = df["destination_center"].apply(lambda x : x[3:9] )
```

```
<ipython-input-651-55eb0ba97d03>:4: FutureWarning: The default value of regex will
change from True to False in a future version. In addition, single character regul
ar expressions will *not* be treated as literal strings when regex=True.
  df["source_state"] = df["source_name"].str.split(" ",n=1,expand=True)[1].str.rep
lace(",").str.replace("", "")
<ipython-input-651-55eb0ba97d03>:7: FutureWarning: The default value of regex will
change from True to False in a future version. In addition, single character regul
ar expressions will *not* be treated as literal strings when regex=True.
  df["destination_state"] = df["destination_name"].str.split(" ",n=1,expand=True)
[1].str.replace(",").str.replace("", "")
```

```
In [ ]: #Time taken between the trip end and start time
```

```
df["time_taken_btwn_odstart_and_od_end"] = ((df["od_end_time"]-df["od_start_time"]))
```

```
In [ ]: #Converting given time duration columns into hours
```

```
df["start_scan_to_end_scan"] = df["start_scan_to_end_scan"]/60
df["actual_time"] = df["actual_time"]/60
df["osrm_time"] = df["osrm_time"]/60
df["segment_actual_time"] = df["segment_actual_time"]/60
df["segment_osrm_time"] = df["segment_osrm_time"]/60
```

```
In [ ]: #Have a glance at modified dataset
```

```
df.head()
```

Out[ ]:

	data	trip_creation_time	route_schedule_uuid	route_type	trip_uuid	source_cei
0	training	2018-09-20 02:35:36.476840	thanos::sroute:eb7bfc78- b351-4c0e-a951- fa3d5c3...	Carting	trip- 153741093647649320	IND388121/
1	training	2018-09-20 02:35:36.476840	thanos::sroute:eb7bfc78- b351-4c0e-a951- fa3d5c3...	Carting	trip- 153741093647649320	IND388121/
2	training	2018-09-20 02:35:36.476840	thanos::sroute:eb7bfc78- b351-4c0e-a951- fa3d5c3...	Carting	trip- 153741093647649320	IND388121/
3	training	2018-09-20 02:35:36.476840	thanos::sroute:eb7bfc78- b351-4c0e-a951- fa3d5c3...	Carting	trip- 153741093647649320	IND388121/
4	training	2018-09-20 02:35:36.476840	thanos::sroute:eb7bfc78- b351-4c0e-a951- fa3d5c3...	Carting	trip- 153741093647649320	IND388121/

5 rows × 31 columns

In [ ]: *#Analysing Dataset after feature creation*

```
df.info()
```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 144867 entries, 0 to 144866
Data columns (total 31 columns):
 #   Column           Non-Null Count Dtype  
--- 
 0   data              144867 non-null  object  
 1   trip_creation_time 144867 non-null  datetime64[ns]
 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  datetime64[ns]
 10  od_end_time         144867 non-null  datetime64[ns]
 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 
 24  source_city           144574 non-null  object  
 25  source_state          144574 non-null  object  
 26  destination_city      144606 non-null  object  
 27  destination_state     144606 non-null  object  
 28  source_pincode        144867 non-null  object  
 29  destination_pincode   144867 non-null  object  
 30  time_taken_btwn_odstart_and_od_end 144867 non-null  float64 
dtypes: bool(1), datetime64[ns](3), float64(11), int64(1), object(15)
memory usage: 33.3+ MB

```

## Data Cleaning

```

In [ ]: #Updating source state values
df["source_state"] = df["source_state"].replace({"Goa Goa":"Goa",
                                                "Layout PC Karnataka":"Karnataka",
                                                "Vadgaon Sheri DPC Maharashtra":"Maharashtra",
                                                "Pashan DPC Maharashtra":"Maharashtra",
                                                "City Madhya Pradesh":"Madhya Pradesh",
                                                "02_DPC Uttar Pradesh":"Uttar Pradesh",
                                                "Nagar_DC Rajasthan":"Rajasthan",
                                                "Alipore_DPC West Bengal":"West Bengal",
                                                "Mandakni Madhya Pradesh":"Madhya Pradesh",
                                                "West _Dc Maharashtra":"Maharashtra",
                                                "DC Rajasthan":"Rajasthan",
                                                "MP Nagar Madhya Pradesh":"Madhya Pradesh",
                                                "Antop Hill Maharashtra":"Maharashtra",
                                                "Avenue_DPC West Bengal":"West Bengal",
                                                "Nagar Uttar Pradesh":"Uttar Pradesh",
                                                "Balaji Nagar Maharashtra":"Maharashtra",
                                                "Kothanur_L Karnataka":"Karnataka",
                                                "Rahatani DPC Maharashtra":"Maharashtra",
                                                "Mahim Maharashtra":"Maharashtra",
                                                "DC Maharashtra":"Maharashtra",
                                                "_NAD Andhra Pradesh":"Andhra Pradesh",
                                               })

```

```
#Updating destination state values
df["destination_state"] = df["destination_state"].replace({
    "Goa Goa": "Goa",
    "Layout PC Karnataka": "Karnataka",
    "Vadgaon Sheri DPC Maharashtra": "Maharashtra",
    "Pashan DPC Maharashtra": "Maharashtra",
    "City Madhya Pradesh": "Madhya Pradesh",
    "02_DPC Uttar Pradesh": "Uttar Pradesh",
    "Nagar_DC Rajasthan": "Rajasthan",
    "Alipore_DPC West Bengal": "West Bengal",
    "Mandakni Madhya Pradesh": "Madhya Pradesh",
    "West _Dc Maharashtra": "Maharashtra",
    "DC Rajasthan": "Rajasthan",
    "MP Nagar Madhya Pradesh": "Madhya Pradesh",
    "Antop Hill Maharashtra": "Maharashtra",
    "Avenue_DPC West Bengal": "West Bengal",
    "Nagar Uttar Pradesh": "Uttar Pradesh",
    "Balaji Nagar Maharashtra": "Maharashtra",
    "Kothanur_L Karnataka": "Karnataka",
    "Rahatani DPC Maharashtra": "Maharashtra",
    "Mahim Maharashtra": "Maharashtra",
    "DC Maharashtra": "Maharashtra",
    "_NAD Andhra Pradesh": "Andhra Pradesh",
    "Delhi Delhi": "Delhi",
    "West_Dc Maharashtra": "Maharashtra",
    "Hub Maharashtra": "Maharashtra"
})
```

In [ ]: #Updating destination and source city values

```
df["destination_city"].replace({
    "del": "Delhi"
}, inplace=True)
df["source_city"].replace({
    "del": "Delhi"
}, inplace=True)

df["source_city"].replace({
    "Bangalore": "Bengaluru"
}, inplace=True)
df["destination_city"].replace({
    "Bangalore": "Bengaluru"
}, inplace=True)
df["destination_city"].replace({
    "AMD": "Ahmedabad"
}, inplace=True)
df["destination_city"].replace({
    "Amdavad": "Ahmedabad"
}, inplace=True)
df["source_city"].replace({
    "AMD": "Ahmedabad"
}, inplace=True)
df["source_city"].replace({
    "Amdavad": "Ahmedabad"
}, inplace=True)
```

In [ ]: # Creating Feature - (Source city + state & Destination city + state)

```
df["source_city_state"] = df["source_city"] + " " + df["source_state"]
df["destination_city_state"] = df["destination_city"] + " " + df["destination_state"]
```

In [ ]: df["source\_city\_state"].nunique()

```
Out[ ]: 1249
```

```
In [ ]: df["destination_city_state"].nunique()
```

```
Out[ ]: 1242
```

```
In [ ]: df["source_state"].nunique()
```

```
Out[ ]: 33
```

```
In [ ]: df["destination_state"].nunique()
```

```
Out[ ]: 32
```

Company delivers in approximately all the states and cities of India.

```
In [ ]: #Dropping Unnecessary columns
```

```
data = df.copy()
```

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

```
Out[ ]: (144867, 33)
```

```
In [ ]: data.drop(  
    ['source_center', "source_name", "destination_center", "destination_name", "cutoff_  
    axis = 1,  
    inplace=True  
)
```

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

```
Out[ ]: (144867, 26)
```

```
In [ ]: #Merging of rows and aggregation of fields
```

```
#Calculating total sum of actual time for each unique trip  
actual_time = data.groupby(["trip_uuid",  
    "start_scan_to_end_scan"])["actual_time"].max().reset_index().groupby  
actual_time
```

Out[ ]:

	trip_uuid	actual_time
0	trip-153671041653548748	26.033333
1	trip-153671042288605164	2.383333
2	trip-153671043369099517	55.783333
3	trip-153671046011330457	0.983333
4	trip-153671052974046625	5.683333
...	...	...
14812	trip-153861095625827784	1.383333
14813	trip-153861104386292051	0.350000
14814	trip-153861106442901555	4.700000
14815	trip-153861115439069069	4.400000
14816	trip-153861118270144424	4.583333

14817 rows × 2 columns

In [ ]:

#Calculating total sum of OSRM segment time for each unique trip

```
segment_osrm_time = data[["trip_uuid","segment_osrm_time"]].groupby("trip_uuid")["segment_osrm_time"]
```

Out[ ]:

	trip_uuid	segment_osrm_time
0	trip-153671041653548748	16.800000
1	trip-153671042288605164	1.083333
2	trip-153671043369099517	32.350000
3	trip-153671046011330457	0.266667
4	trip-153671052974046625	1.916667
...	...	...
14812	trip-153861095625827784	1.033333
14813	trip-153861104386292051	0.183333
14814	trip-153861106442901555	1.466667
14815	trip-153861115439069069	3.683333
14816	trip-153861118270144424	1.116667

14817 rows × 2 columns

In [ ]:

#Calculating total sum of actual segment time for each unique trip

```
segment_actual_time = data.groupby("trip_uuid")["segment_actual_time"].sum().reset_index()
```

Out[ ]:

	trip_uuid	segment_actual_time
0	trip-153671041653548748	25.800000
1	trip-153671042288605164	2.350000
2	trip-153671043369099517	55.133333
3	trip-153671046011330457	0.983333
4	trip-153671052974046625	5.666667
...	...	...
14812	trip-153861095625827784	1.366667
14813	trip-153861104386292051	0.350000
14814	trip-153861106442901555	4.683333
14815	trip-153861115439069069	4.300000
14816	trip-153861118270144424	4.566667

14817 rows × 2 columns

In [ ]:

```
#Calculating total sum of OSRM time for each unique trip

osrm_time = data.groupby(["trip_uuid",
                           "start_scan_to_end_scan"])["osrm_time"].max().reset_index().groupby('
osrm_time
```

Out[ ]:

	trip_uuid	osrm_time
0	trip-153671041653548748	12.383333
1	trip-153671042288605164	1.133333
2	trip-153671043369099517	29.016667
3	trip-153671046011330457	0.250000
4	trip-153671052974046625	1.950000
...	...	...
14812	trip-153861095625827784	1.033333
14813	trip-153861104386292051	0.200000
14814	trip-153861106442901555	0.900000
14815	trip-153861115439069069	3.066667
14816	trip-153861118270144424	1.133333

14817 rows × 2 columns

In [ ]:

```
time_taken_btwn_odstart_and_od_end = data.groupby("trip_uuid")["time_taken_btwn_ods
time_taken_btwn_odstart_and_od_end
```

```
Out[ ]:
```

	trip_uuid	time_taken_btwn_odstart_and_od_end
0	trip-153671041653548748	[16.65842298, 21.0100736875]
1	trip-153671042288605164	[2.0463247669444447, 0.9805397955555556]
2	trip-153671043369099517	[51.662059856388886, 13.910648811388889]
3	trip-153671046011330457	[1.6749155866666667]
4	trip-153671052974046625	[2.5335485744444446, 1.342388563333332, 8.096...]
...	...	...
14812	trip-153861095625827784	[2.546464057777778, 1.7540180775]
14813	trip-153861104386292051	[1.0098420219444444]
14814	trip-153861106442901555	[2.895179575833333, 4.1401515375]
14815	trip-153861115439069069	[1.7609491794444445, 0.7362400538888889, 1.035...]
14816	trip-153861118270144424	[1.1155594141666667, 4.7912334425]

14817 rows × 2 columns

```
In [ ]: time_taken_btwn_odstart_and_od_end["time_taken_btwn_odstart_and_od_end"] = time_tak  
time_taken_btwn_odstart_and_od_end["time_taken_btwn_odstart_and_od_end"]
```

```
Out[ ]: 0      37.668497  
1      3.026865  
2      65.572709  
3      1.674916  
4      11.972484  
      ...  
14812     4.300482  
14813     1.009842  
14814     7.035331  
14815     5.808548  
14816     5.906793  
Name: time_taken_btwn_odstart_and_od_end, Length: 14817, dtype: float64
```

```
In [ ]: start_scan_to_end_scan = ((data.groupby("trip_uuid")["start_scan_to_end_scan"].unique()  
start_scan_to_end_scan
```

Out[ ]:

	trip_uuid	start_scan_to_end_scan
0	trip-153671041653548748	[16.65, 21.0]
1	trip-153671042288605164	[2.033333333333333, 0.9666666666666667]
2	trip-153671043369099517	[51.65, 13.9]
3	trip-153671046011330457	[1.6666666666666667]
4	trip-153671052974046625	[2.533333333333333, 1.333333333333333, 8.0833...]
...	...	...
14812	trip-153861095625827784	[2.533333333333333, 1.75]
14813	trip-153861104386292051	[1.0]
14814	trip-153861106442901555	[2.883333333333333, 4.13333333333334]
14815	trip-153861115439069069	[1.75, 0.733333333333333, 1.03333333333334,...]
14816	trip-153861118270144424	[1.1, 4.78333333333333]

14817 rows × 2 columns

In [ ]: `start_scan_to_end_scan["start_scan_to_end_scan"] = start_scan_to_end_scan["start_scan_to_end_scan"]`

Out[ ]: 0 37.650000  
1 3.000000  
2 65.550000  
3 1.666667  
4 11.950000  
...  
14812 4.283333  
14813 1.000000  
14814 7.016667  
14815 5.783333  
14816 5.883333  
Name: start\_scan\_to\_end\_scan, Length: 14817, dtype: float64

In [ ]: `#Calculating total sum of OSRM distance for each unique trip`

```
osrm_distance = data.groupby(["trip_uuid",
                               "start_scan_to_end_scan"])["osrm_distance"].max().reset_index().groupby("trip_uuid").sum()
```

Out[ ]:

	trip_uuid	osrm_distance
0	trip-153671041653548748	991.3523
1	trip-153671042288605164	85.1110
2	trip-153671043369099517	2372.0852
3	trip-153671046011330457	19.6800
4	trip-153671052974046625	146.7918
...	...	...
14812	trip-153861095625827784	73.4630
14813	trip-153861104386292051	16.0882
14814	trip-153861106442901555	63.2841
14815	trip-153861115439069069	177.6635
14816	trip-153861118270144424	80.5787

14817 rows × 2 columns

In [ ]:

```
#Calculating total sum of actual distance for each unique trip

actual_distance_to_destination = data.groupby(["trip_uuid",
                                                "start_scan_to_end_scan"])["actual_distance_to_destination"].max().reset_index()

actual_distance_to_destination
```

Out[ ]:

	trip_uuid	actual_distance_to_destination
0	trip-153671041653548748	824.732854
1	trip-153671042288605164	73.186911
2	trip-153671043369099517	1932.273969
3	trip-153671046011330457	17.175274
4	trip-153671052974046625	127.448500
...	...	...
14812	trip-153861095625827784	57.762332
14813	trip-153861104386292051	15.513784
14814	trip-153861106442901555	38.684839
14815	trip-153861115439069069	134.723836
14816	trip-153861118270144424	66.081533

14817 rows × 2 columns

In [ ]:

```
#Calculating total sum of segment OSRM distance for each unique trip
```

```
segment_osrm_distance = data[["trip_uuid",
                               "segment_osrm_distance"]].groupby("trip_uuid")["segme
```

```
segment_osrm_distance
```

Out[ ]:

	trip_uuid	segment_osrm_distance
0	trip-153671041653548748	1320.4733
1	trip-153671042288605164	84.1894
2	trip-153671043369099517	2545.2678
3	trip-153671046011330457	19.8766
4	trip-153671052974046625	146.7919
...	...	...
14812	trip-153861095625827784	64.8551
14813	trip-153861104386292051	16.0883
14814	trip-153861106442901555	104.8866
14815	trip-153861115439069069	223.5324
14816	trip-153861118270144424	80.5787

14817 rows × 2 columns

## Hypothesis Testing

### 1. Analysing TimeTaken Between OdStart and OdEnd time & StartScanToEndScan :

H0: Mean of time taken between trip end and start time = Mean of start and end scan time

Ha: Mean of time taken between trip end and start time != Mean of start and end scan time

In [ ]: #Visualization for two features

```
plt.figure(figsize=(15,5))
plt.subplot(121)
sns.distplot((time_taken_btwn_odstart_and_od_end["time_taken_btwn_odstart_and_od_end"]))
plt.subplot(122)
sns.distplot((start_scan_to_end_scan["start_scan_to_end_scan"]))

plt.show()
```

```

<ipython-input-678-0daec2699fca>:5: UserWarning:
`distplot` is a deprecated function and will be removed in seaborn v0.14.0.
Please adapt your code to use either `displot` (a figure-level function with
similar flexibility) or `histplot` (an axes-level function for histograms).

For a guide to updating your code to use the new functions, please see
https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751

    sns.distplot((time_taken_btwn_odstart_and_od_end["time_taken_btwn_odstart_and_od_end"]))
<ipython-input-678-0daec2699fca>:7: UserWarning:
`distplot` is a deprecated function and will be removed in seaborn v0.14.0.
Please adapt your code to use either `displot` (a figure-level function with
similar flexibility) or `histplot` (an axes-level function for histograms).

For a guide to updating your code to use the new functions, please see
https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751

    sns.distplot((start_scan_to_end_scan["start_scan_to_end_scan"]))

```

In [ ]: #Performing a Kolmogorov-Smirnov (KS) test

```

from scipy import stats
ks_test, p_value = stats.ks_2samp(time_taken_btwn_odstart_and_od_end["time_taken_bt",
                                                                start_scan_to_end_scan["start_scan_to_end_scan"]])

# Ho: The distribution are similar
# Ha: The distributions are different

if p_value < 0.05:
    print("Reject Ho: The distribution are different.")
else :
    print("Fail to reject Ho: The distribution is same.")

```

Fail to reject Ho: The distribution is same.

In [ ]: #Performing two sample ttest 5 times for 3000 random samples

```

for i in range(5):
    print(stats.ttest_ind((time_taken_btwn_odstart_and_od_end["time_taken_btwn_odst",
                                                               (start_scan_to_end_scan["start_scan_to_end_scan"].sample(3000))]))

```

```

TtestResult(statistic=0.154223020834327, pvalue=0.8774390800218714, df=5998.0)
TtestResult(statistic=-1.6978829946133276, pvalue=0.08958171916692721, df=5998.0)
TtestResult(statistic=-0.541370490611641, pvalue=0.588272349399451, df=5998.0)
TtestResult(statistic=0.8246184631642366, pvalue=0.40962102728226235, df=5998.0)
TtestResult(statistic=0.054966205032367293, pvalue=0.9561672182374987, df=5998.0)

```

We can conclude that there is no significant difference between the average time taken between OD start and end and the average start scan to end scan durations for the population, based on the two-sample t-test.

```
In [ ]: #Checking mean and variance for both features  
time_taken_btwn_odstart_and_od_end["time_taken_btwn_odstart_and_od_end"].mean(),tim  
Out[ ]: (8.861857235305113, 10.981665759990634)  
  
In [ ]: start_scan_to_end_scan["start_scan_to_end_scan"].mean(),start_scan_to_end_scan["sta  
Out[ ]: (8.835777597804324, 10.976286391439693)
```

We observe that the variances and means of both the scan time and the time taken for trip start and end are closely similar.

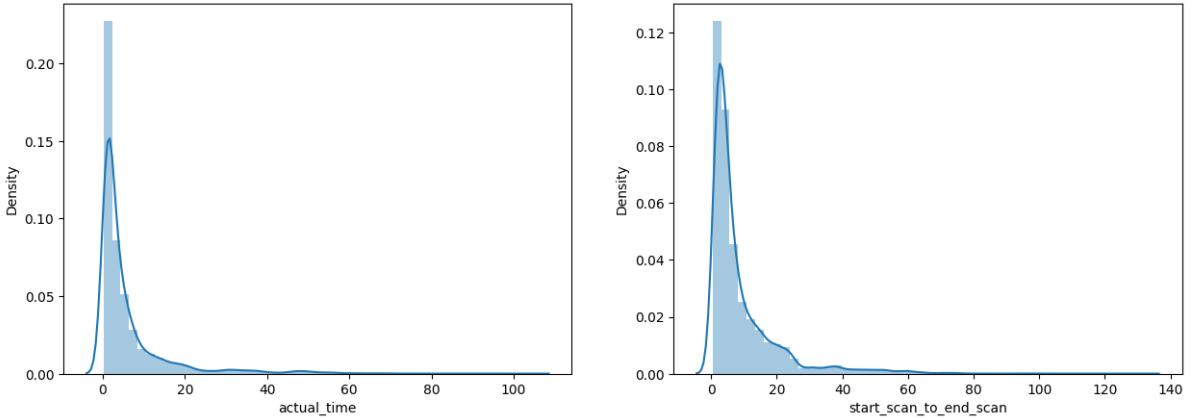
## 2. Analysing Actual Time taken to complete the delivery & start-scan-end-scan

H0: Mean of start and end scan time <= Mean of Actual time taken to complete delivery

Ha: Mean of start and end scan time > Mean of Actual time taken to complete delivery

```
In [ ]: #Visualization for two features  
plt.figure(figsize=(15,5))  
plt.subplot(121)  
sns.distplot((actual_time["actual_time"]))  
plt.subplot(122)  
sns.distplot((start_scan_to_end_scan["start_scan_to_end_scan"]))  
plt.show()
```

```
<ipython-input-683-379510bae572>:5: UserWarning:  
`distplot` is a deprecated function and will be removed in seaborn v0.14.0.  
Please adapt your code to use either `displot` (a figure-level function with  
similar flexibility) or `histplot` (an axes-level function for histograms).  
For a guide to updating your code to use the new functions, please see  
https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751  
sns.distplot((actual_time["actual_time"]))  
<ipython-input-683-379510bae572>:7: UserWarning:  
`distplot` is a deprecated function and will be removed in seaborn v0.14.0.  
Please adapt your code to use either `displot` (a figure-level function with  
similar flexibility) or `histplot` (an axes-level function for histograms).  
For a guide to updating your code to use the new functions, please see  
https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751  
sns.distplot((start_scan_to_end_scan["start_scan_to_end_scan"]))
```



```
In [ ]: #Performing a Kolmogorov-Smirnov (KS) test
stats.ks_2samp(actual_time["actual_time"],start_scan_to_end_scan["start_scan_to_end_scan"])
# Ho: The distribution are similar
# Ha: The distributions are different
if p_value < 0.05:
    print("Reject Ho: The distribution are different.")
else :
    print("Fail to reject Ho: The distribution is same.")
Fail to reject Ho: The distribution is same.
```

```
In [ ]: #Performing two sample ttest 5 times for 3000 random samples
for i in range(5):
    print(stats.ttest_ind((actual_time["actual_time"].sample(3000))
                           ,(start_scan_to_end_scan["start_scan_to_end_scan"].sample(3000)),alternative='less'))
TtestResult(statistic=-10.32493175352972, pvalue=4.378724263421445e-25, df=5998.0)
TtestResult(statistic=-10.398024474271098, pvalue=2.066635659571137e-25, df=5998.0)
TtestResult(statistic=-10.946118176259318, pvalue=6.321168505493412e-28, df=5998.0)
TtestResult(statistic=-10.02224426860614, pvalue=9.3002614704782e-24, df=5998.0)
TtestResult(statistic=-11.308329092650068, pvalue=1.180848508301586e-29, df=5998.0)
```

```
In [ ]: #Checking mean and variance for both features
actual_time["actual_time"].mean(),actual_time["actual_time"].std()
Out[ ]: (5.945176711435065, 9.355547822973838)
```

```
In [ ]: start_scan_to_end_scan["start_scan_to_end_scan"].mean(),start_scan_to_end_scan["start_scan_to_end_scan"].std()
Out[ ]: (8.835777597804324, 10.976286391439693)
```

We can conclude that, on average, the "actual\_time" is significantly less than the "start\_scan\_to\_end\_scan".

### 3. Analysing Actual Time & Time Taken between start and end trip time.

H0: Mean of Actual time taken to complete delivery = Mean of time taken between trip end and start time

Ha: Mean of Actual time taken to complete delivery != Mean of time taken between trip end and start time

In [ ]: #Visualization for two features

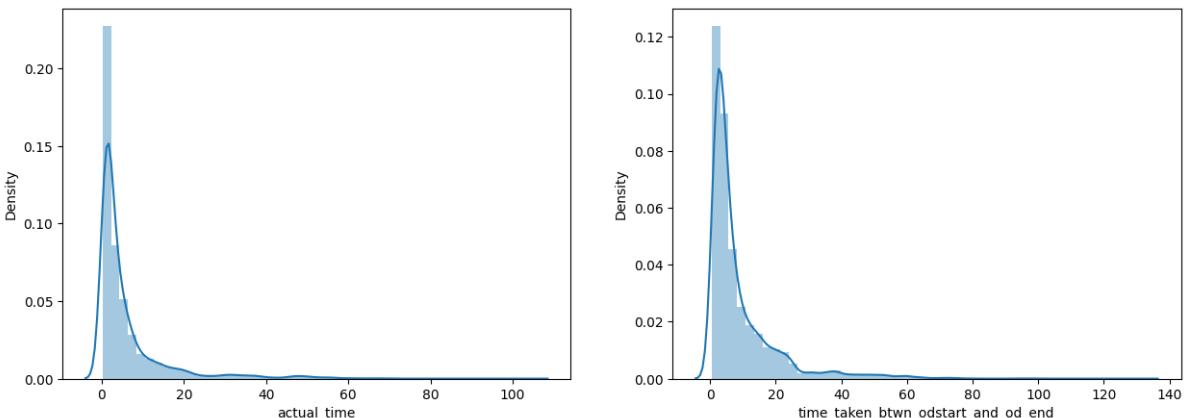
```
plt.figure(figsize=(15,5))
plt.subplot(121)
sns.distplot((actual_time["actual_time"]))
plt.subplot(122)
sns.distplot((time_taken_btwn_odstart_and_od_end["time_taken_btwn_odstart_and_od_end"]))

plt.show()
```

```
<ipython-input-688-22f0bf166288>:5: UserWarning:
`distplot` is a deprecated function and will be removed in seaborn v0.14.0.
Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
For a guide to updating your code to use the new functions, please see
https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751

sns.distplot((actual_time["actual_time"]))
<ipython-input-688-22f0bf166288>:7: UserWarning:
`distplot` is a deprecated function and will be removed in seaborn v0.14.0.
Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
For a guide to updating your code to use the new functions, please see
https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751

sns.distplot((time_taken_btwn_odstart_and_od_end["time_taken_btwn_odstart_and_od_end"]))
```



In [ ]: #Performing a Kolmogorov-Smirnov (KS) test

```
stats.ks_2samp(actual_time["actual_time"],time_taken_btwn_odstart_and_od_end["time_taken_btwn_odstart_and_od_end"])

# Ho: The distribution are similar
# Ha: The distributions are different

if p_value < 0.05:
    print("Reject Ho: The distribution are different.")
else :
    print("Fail to reject Ho: The distribution is same.")
```

Fail to reject H<sub>0</sub>: The distribution is same.

```
In [ ]: #Performing two sample ttest 5 times for 3000 random samples

for i in range(5):
    print(stats.ttest_ind((actual_time["actual_time"].sample(3000)),
                          (time_taken_btwn_odstart_and_od_end["time_taken_btwn_odstart_and_o
TtestResult(statistic=-12.518538378705012, pvalue=1.6377161813697582e-35, df=5998.
0)
TtestResult(statistic=-11.843110488298858, pvalue=5.298735991249866e-32, df=5998.
0)
TtestResult(statistic=-9.911565571121468, pvalue=5.564830735171225e-23, df=5998.0)
TtestResult(statistic=-11.32087890166514, pvalue=2.0529595074053258e-29, df=5998.
0)
TtestResult(statistic=-12.101002481954243, pvalue=2.5443444293753415e-33, df=5998.
0)
```

we can conclude that there is a statistically significant difference between the average "actual\_time" and the average "time\_taken\_btwn\_odstart\_and\_od\_end".

#### 4. Analysing Actual Time taken to complete delivery from source to destination hub & OSRM measured time :

H<sub>0</sub>: Mean of OSRM time  $\geq$  Mean of Actual time taken to complete delivery

H<sub>a</sub>: Mean of OSRM time  $<$  Mean of Actual time taken to complete delivery

```
In [ ]: #Visulization for two features
```

```
plt.figure(figsize=(10,4))
plt.subplot(121)
sns.distplot(((actual_time["actual_time"])))
plt.subplot(122)
sns.distplot(((osrm_time["osrm_time"])))

plt.show()
```

```
<ipython-input-691-9f52ded87e48>:5: UserWarning:
```

```
`distplot` is a deprecated function and will be removed in seaborn v0.14.0.
```

```
Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
```

```
For a guide to updating your code to use the new functions, please see  
https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751
```

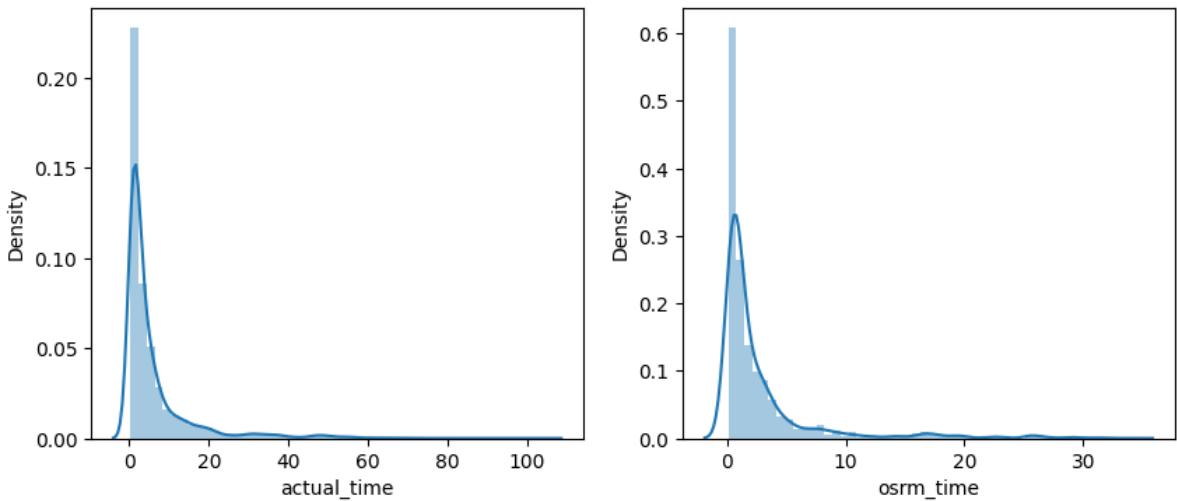
```
    sns.distplot(((actual_time["actual_time"])))
<ipython-input-691-9f52ded87e48>:7: UserWarning:
```

```
`distplot` is a deprecated function and will be removed in seaborn v0.14.0.
```

```
Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
```

```
For a guide to updating your code to use the new functions, please see  
https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751
```

```
    sns.distplot(((osrm_time["osrm_time"])))
```



```
In [ ]: #Performing a Kolmogorov-Smirnov (KS) test
stats.ks_2samp(actual_time["actual_time"],osrm_time["osrm_time"])

# Ho: The distribution are similar
# Ha: The disbutions are different

if p_value < 0.05:
    print("Reject Ho: The distribution are different.")
else :
    print("Fail to reject Ho: The distribution is same.")

Fail to reject Ho: The distribution is same.
```

```
In [ ]: #Performing two sample ttest 5 times for 3000 random samples
for i in range(5):
    print(stats.ttest_ind(actual_time["actual_time"].sample(3000),
                          osrm_time["osrm_time"].sample(3000),alternative='greater'))

TtestResult(statistic=19.331026256687522, pvalue=4.060844408979209e-81, df=5998.0)
TtestResult(statistic=17.679150302545757, pvalue=1.5936959951365007e-68, df=5998.0)
TtestResult(statistic=17.65903160800764, pvalue=2.236747260467911e-68, df=5998.0)
TtestResult(statistic=17.717982889391347, pvalue=8.276329241068681e-69, df=5998.0)
TtestResult(statistic=17.898970470259492, pvalue=3.8387546445626664e-70, df=5998.0)
```

We can conclude that, on average, the "actual\_time" is significantly greater than the "osrm\_time".

```
In [ ]: #Checking mean and variance for both features
actual_time["actual_time"].mean(),actual_time["actual_time"].std()
```

```
Out[ ]: (5.945176711435065, 9.355547822973838)
```

```
In [ ]: osrm_time["osrm_time"].mean(),osrm_time["osrm_time"].std()
Out[ ]: (2.6973138962003107, 4.537654251845704)
```

## 5. Analysing Actual Time taken to complete delivery from source to destination hub & Segment Actual Time :

H0: Actual time = segment actual time

Ha: Actual time != segment actual time

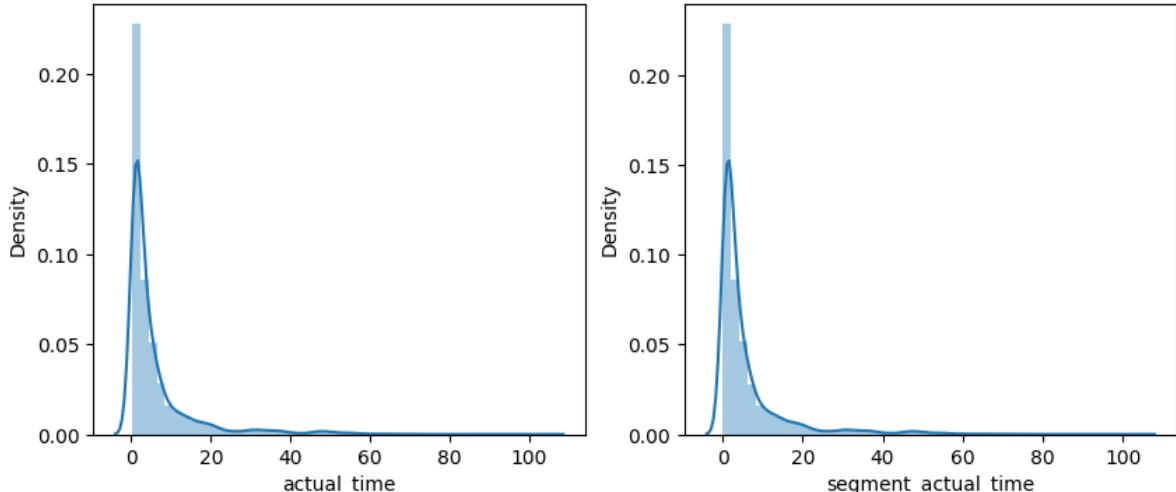
```
In [ ]: plt.figure(figsize=(10,4))
plt.subplot(121)
sns.distplot(((actual_time["actual_time"])))
plt.subplot(122)
sns.distplot(((segment_actual_time["segment_actual_time"])))

plt.show()
```

```
<ipython-input-696-1fcfd05e8e5d>:3: UserWarning:
`distplot` is a deprecated function and will be removed in seaborn v0.14.0.
Please adapt your code to use either `displot` (a figure-level function with
similar flexibility) or `histplot` (an axes-level function for histograms).
For a guide to updating your code to use the new functions, please see
https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751

sns.distplot(((actual_time["actual_time"])))
<ipython-input-696-1fcfd05e8e5d>:5: UserWarning:
`distplot` is a deprecated function and will be removed in seaborn v0.14.0.
Please adapt your code to use either `displot` (a figure-level function with
similar flexibility) or `histplot` (an axes-level function for histograms).
For a guide to updating your code to use the new functions, please see
https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751

sns.distplot(((segment_actual_time["segment_actual_time"])))
```



```
In [ ]: for i in range(7):
    print(stats.ttest_ind((actual_time["actual_time"].sample(3000)),
                          (segment_actual_time["segment_actual_time"].sample(3000))))
```

```
TtestResult(statistic=-0.04403168012141797, pvalue=0.964880616979656, df=5998.0)
TtestResult(statistic=1.197278661844238, pvalue=0.23124527379569798, df=5998.0)
TtestResult(statistic=0.7478989135948552, pvalue=0.45455045430691066, df=5998.0)
TtestResult(statistic=0.954723389172665, pvalue=0.3397560735882553, df=5998.0)
TtestResult(statistic=0.4138476003794712, pvalue=0.6790005179083995, df=5998.0)
TtestResult(statistic=-0.9863135871935026, pvalue=0.32401903874986027, df=5998.0)
TtestResult(statistic=-0.5560487664432213, pvalue=0.5781982590967985, df=5998.0)
```

From two sample ttest , we can conclude that Population average for Actual Time taken to complete delivery trip and segment actual time are same.

```
In [ ]: actual_time["actual_time"].mean(),actual_time["actual_time"].std()
```

```
Out[ ]: (5.945176711435065, 9.355547822973838)
```

```
In [ ]: segment_actual_time["segment_actual_time"].mean(),segment_actual_time["segment_actu
```

```
Out[ ]: (5.8982047647971925, 9.27079941315281)
```

## 6. Analysing osrm Time & segment-osrm-time :

Ho: segment actual time  $\leq$  OSRM time

Ha: segment actual time  $>$  OSRM time

```
In [ ]: plt.figure(figsize=(10,4))
plt.subplot(121)
sns.distplot(((osrm_time["osrm_time"])))
plt.subplot(122)
sns.distplot(((segment_osrm_time["segment_osrm_time"])))

plt.show()
```

```
<ipython-input-700-cdfa1cbadd0>:3: UserWarning:
```

```
`distplot` is a deprecated function and will be removed in seaborn v0.14.0.
```

```
Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
```

```
For a guide to updating your code to use the new functions, please see  
https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751
```

```
    sns.distplot(((osrm_time["osrm_time"])))
```

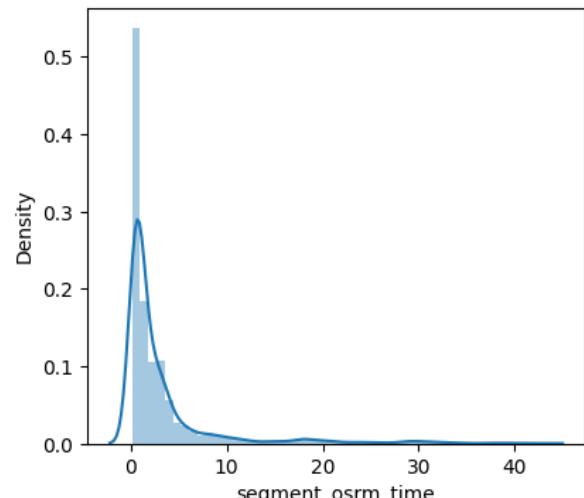
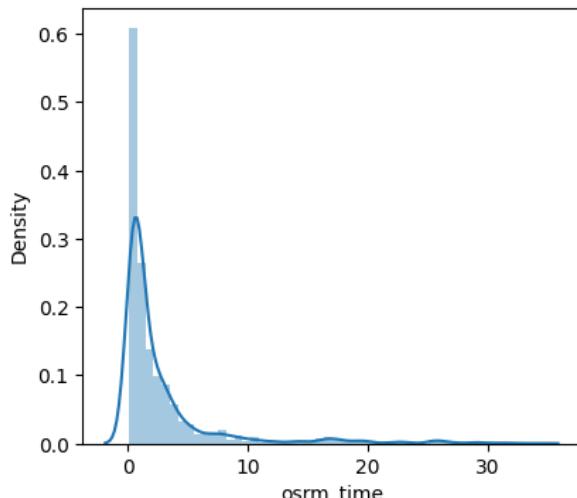
```
<ipython-input-700-cdfa1cbadd0>:5: UserWarning:
```

```
`distplot` is a deprecated function and will be removed in seaborn v0.14.0.
```

```
Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
```

```
For a guide to updating your code to use the new functions, please see  
https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751
```

```
    sns.distplot(((segment_osrm_time["segment_osrm_time"])))
```



```
In [ ]: for i in range(7):
    print(stats.ttest_ind((osrm_time["osrm_time"].sample(3000)),
                          (segment_osrm_time["segment_osrm_time"].sample(3000)), alternative =
TtestResult(statistic=-3.090642400486092, pvalue=0.0010031969814963257, df=5998.0)
TtestResult(statistic=-2.359448729941225, pvalue=0.009166987446761378, df=5998.0)
TtestResult(statistic=-3.0613692444184664, pvalue=0.0011065118011126605, df=5998.0)
TtestResult(statistic=-0.5219608861811813, pvalue=0.30085841925025064, df=5998.0)
TtestResult(statistic=-1.2762716273895738, pvalue=0.10095446577214666, df=5998.0)
TtestResult(statistic=-1.8788767304515988, pvalue=0.030154890703798506, df=5998.0)
TtestResult(statistic=-2.8940803408564344, pvalue=0.0019082090073203265, df=5998.0)
```

```
In [ ]: osrm_time["osrm_time"].mean(),osrm_time["osrm_time"].std()
```

```
Out[ ]: (2.6973138962003107, 4.537654251845704)
```

```
In [ ]: segment_osrm_time["segment_osrm_time"].mean(),segment_osrm_time["segment_osrm_time"]
```

```
Out[ ]: (3.0158297901059594, 5.242367441693001)
```

From ttest , we can conclude that average of osrm Time & segment-osrm-time for population is not same.

## 7. Analysing and Visulizing OSRM Estimated distance and Segment-osrm-distance :

H0 : Segment OSRM distnace  $\leq$  OSRM distnace

H<sub>a</sub> : Segment OSRM distnace  $>$  OSRM distnace

```
In [ ]: plt.figure(figsize=(10,4))
plt.subplot(121)
sns.distplot(((osrm_distance["osrm_distance"])))
plt.subplot(122)
sns.distplot(((segment_osrm_distance["segment_osrm_distance"])))

plt.show()
```

```
<ipython-input-704-b59a614a5fb>:3: UserWarning:
```

```
`distplot` is a deprecated function and will be removed in seaborn v0.14.0.
```

```
Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
```

```
For a guide to updating your code to use the new functions, please see  
https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751
```

```
sns.distplot(((osrm_distance["osrm_distance"])))
```

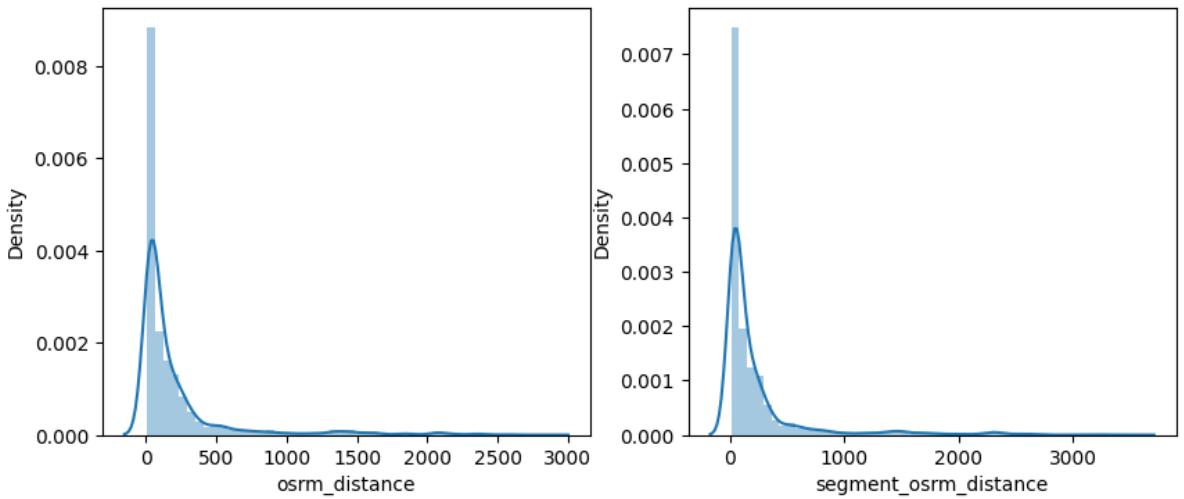
```
<ipython-input-704-b59a614a5fb>:5: UserWarning:
```

```
`distplot` is a deprecated function and will be removed in seaborn v0.14.0.
```

```
Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
```

```
For a guide to updating your code to use the new functions, please see  
https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751
```

```
sns.distplot(((segment_osrm_distance["segment_osrm_distance"])))
```



```
In [ ]: stats.ks_2samp(osrm_distance["osrm_distance"], segment_osrm_distance["segment_osrm_distance"])
Out[ ]: KtestResult(statistic=0.03948167645272321, pvalue=1.8042208791084262e-10, statistic_location=50.2941, statistic_sign=1)
```

```
In [ ]: for i in range(7):
    print(stats.ttest_ind(osrm_distance["osrm_distance"].sample(5000),
                          segment_osrm_distance["segment_osrm_distance"].sample(5000), alternative='less'))
```

TtestResult(statistic=-1.4554027224744035, pvalue=0.0727945999128067, df=9998.0)  
TtestResult(statistic=-3.4819450672215173, pvalue=0.0002499563783931604, df=9998.0)  
TtestResult(statistic=-4.025497420856429, pvalue=2.8637275778193898e-05, df=9998.0)  
TtestResult(statistic=-2.1368881623737823, pvalue=0.01631564393274573, df=9998.0)  
TtestResult(statistic=-2.0232190567426787, pvalue=0.021538562054093827, df=9998.0)  
TtestResult(statistic=-2.7422297315928144, pvalue=0.0030566115332797956, df=9998.0)  
TtestResult(statistic=-1.766707015708861, pvalue=0.03865390712052963, df=9998.0)

```
In [ ]: osrm_distance["osrm_distance"].mean(), osrm_distance["osrm_distance"].std()
Out[ ]: (204.83672531551593, 370.74927471335474)
```

```
In [ ]: segment_osrm_distance["segment_osrm_distance"].mean(), segment_osrm_distance["segment_osrm_distance"].std()
Out[ ]: (223.20116128771005, 416.62837429074216)
```

We can conclude that the distributions of "segment osrm distance" and "osrm distance" are not the same.

We can conclude that the average "osrm distance" for the population is significantly less than the average "segment osrm distance".

## 8. Analysing and Visualizing OSRM Estimated distance and Actual Distance between source and destination warehouse :

H0 : Mean OSRM distance  $\leq$  Mean Actual distance

Ha : Mean OSRM distance  $>$  Mean Actual distance

```
In [ ]: plt.figure(figsize=(10,4))
plt.subplot(121)
sns.distplot(((osrm_distance["osrm_distance"])))
```

```

plt.subplot(122)
sns.distplot(((actual_distance_to_destination["actual_distance_to_destination"])))

plt.show()

```

<ipython-input-709-f43d32a85a66>:3: UserWarning:

`distplot` is a deprecated function and will be removed in seaborn v0.14.0.

Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

For a guide to updating your code to use the new functions, please see <https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751>

```

sns.distplot(((osrm_distance["osrm_distance"])))

<ipython-input-709-f43d32a85a66>:5: UserWarning:
```

`distplot` is a deprecated function and will be removed in seaborn v0.14.0.

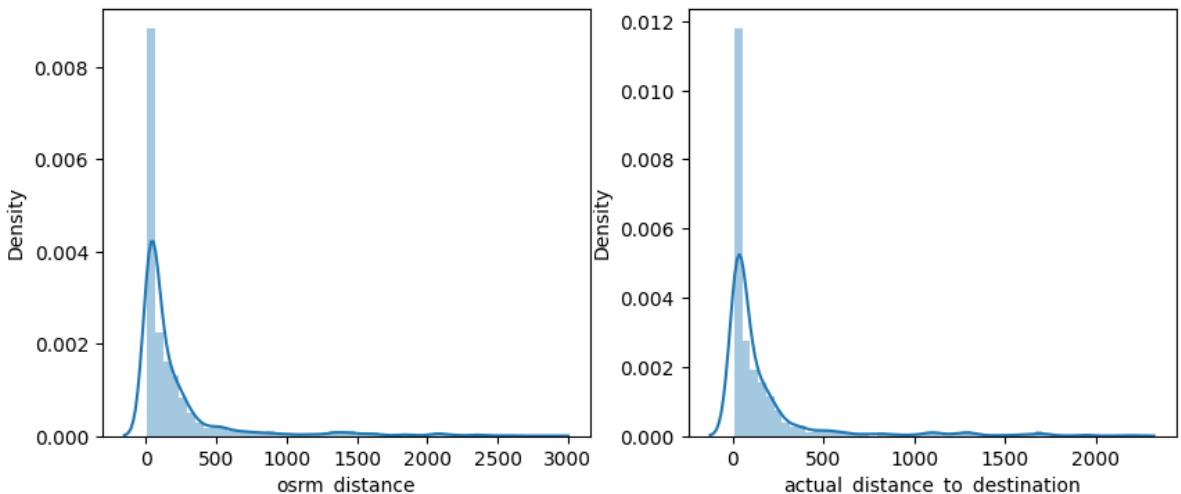
Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

For a guide to updating your code to use the new functions, please see <https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751>

```

sns.distplot(((actual_distance_to_destination["actual_distance_to_destination"])))

```



In [ ]: stats.ks\_2samp(osrm\_distance["osrm\_distance"], actual\_distance\_to\_destination["actual\_distance\_to\_destination"])

Out[ ]: KstestResult(statistic=0.11837753931295136, pvalue=6.578385372142345e-91, statistic\_c(location=25.247485296255537, statistic\_sign=-1)

In [ ]: for i in range(5):
 print(stats.ttest\_ind(osrm\_distance["osrm\_distance"].sample(5000),
 actual\_distance\_to\_destination["actual\_distance\_to\_destination"].sample(5000)))

TtestResult(statistic=6.2300387600042315, pvalue=2.425513267983247e-10, df=9998.0)
TtestResult(statistic=3.536515289190613, pvalue=0.0002036395371832916, df=9998.0)
TtestResult(statistic=6.392236882191162, pvalue=8.538535214416912e-11, df=9998.0)
TtestResult(statistic=5.825796803842123, pvalue=2.929799568567547e-09, df=9998.0)
TtestResult(statistic=7.21911048743959, pvalue=2.8068710698196955e-13, df=9998.0)

In [ ]: osrm\_distance["osrm\_distance"].mean(), osrm\_distance["osrm\_distance"].std()

Out[ ]: (204.83672531551593, 370.74927471335474)

```
In [ ]: actual_distance_to_destination["actual_distance_to_destination"].mean(),actual_dist  
Out[ ]: (164.47332174544243, 305.5408288910492)
```

From left sided ttest , we can conclude for population OSRM estimated distance is higher than the actual distance from source to destination warehouse.

## Merging

```
In [ ]: distances = segment_osrm_distance.merge(actual_distance_to_destination.merge(osrm_c  
on="trip_uuid",  
left_on="trip_uuid",  
right_on="trip_uuid"))
```

```
In [ ]: distances
```

```
Out[ ]:
```

	trip_uuid	segment_osrm_distance	actual_distance_to_destination	osrm_distance
0	trip-153671041653548748	1320.4733	824.732854	991.3523
1	trip-153671042288605164	84.1894	73.186911	85.1110
2	trip-153671043369099517	2545.2678	1932.273969	2372.0852
3	trip-153671046011330457	19.8766	17.175274	19.6800
4	trip-153671052974046625	146.7919	127.448500	146.7918
...	...	...	...	...
14812	trip-153861095625827784	64.8551	57.762332	73.4630
14813	trip-153861104386292051	16.0883	15.513784	16.0882
14814	trip-153861106442901555	104.8866	38.684839	63.2841
14815	trip-153861115439069069	223.5324	134.723836	177.6635
14816	trip-153861118270144424	80.5787	66.081533	80.5787

14817 rows × 4 columns

```
In [ ]: time = segment_osrm_time.merge(osrm_time.merge(segment_actual_time.merge(actual_time  
on="trip_uuid",  
left_on="trip_uuid",  
right_on="trip_uuid"),on="trip_uuid"),on="trip_uuid"),on="trip  
time
```

Out[ ]:

		trip_uuid	segment_osrm_time	osrm_time	segment_actual_time	actual_time	tin	
0		trip-153671041653548748		16.800000	12.383333		25.800000	26.033333
1		trip-153671042288605164		1.083333	1.133333		2.350000	2.383333
2		trip-153671043369099517		32.350000	29.016667		55.133333	55.783333
3		trip-153671046011330457		0.266667	0.250000		0.983333	0.983333
4		trip-153671052974046625		1.916667	1.950000		5.666667	5.683333
	...	...	...	...	...	...	...	...
14812		trip-153861095625827784		1.033333	1.033333		1.366667	1.383333
14813		trip-153861104386292051		0.183333	0.200000		0.350000	0.350000
14814		trip-153861106442901555		1.466667	0.900000		4.683333	4.700000
14815		trip-153861115439069069		3.683333	3.066667		4.300000	4.400000
14816		trip-153861118270144424		1.116667	1.133333		4.566667	4.583333

14817 rows × 7 columns

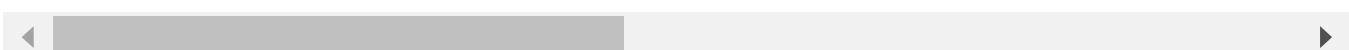


In [ ]: `Merge1 = time.merge(distances, on="trip_uuid", )  
Merge1`

Out[ ]:

		trip_uuid	segment_osrm_time	osrm_time	segment_actual_time	actual_time	tin
0		trip-153671041653548748		16.800000	12.383333	25.800000	26.033333
1		trip-153671042288605164		1.083333	1.133333	2.350000	2.383333
2		trip-153671043369099517		32.350000	29.016667	55.133333	55.783333
3		trip-153671046011330457		0.266667	0.250000	0.983333	0.983333
4		trip-153671052974046625		1.916667	1.950000	5.666667	5.683333
...		...		...	...	...	...
14812		trip-153861095625827784		1.033333	1.033333	1.366667	1.383333
14813		trip-153861104386292051		0.183333	0.200000	0.350000	0.350000
14814		trip-153861106442901555		1.466667	0.900000	4.683333	4.700000
14815		trip-153861115439069069		3.683333	3.066667	4.300000	4.400000
14816		trip-153861118270144424		1.116667	1.133333	4.566667	4.583333

14817 rows × 10 columns



Merging Location details and route\_type and Numerical data on TripID :

```
In [ ]: city = data.groupby("trip_uuid")[["source_city",
                                         "destination_city"]].aggregate({
    "source_city":pd.unique,
    "destination_city":pd.unique,
})

state = data.groupby("trip_uuid")[["source_state",
                                    "destination_state"]].aggregate({
    "source_state":pd.unique,
    "destination_state":pd.unique,
})

city_state = data.groupby("trip_uuid")[["source_city_state",
                                         "destination_city_state"]].aggregate({
    "source_city_state":pd.unique,
    "destination_city_state":pd.unique,
})

locations = city.merge(city_state.merge(state,on="trip_uuid",
                                         ,how="outer"),
                       on="trip_uuid",
                       how="outer")
```

```
In [ ]: route_type = data.groupby("trip_uuid")["route_type"].unique().reset_index()

Merged = route_type.merge(locations.merge(Merge1,on="trip_uuid",
```

```
        how="outer"),
        on="trip_uuid",
how="outer"
)
```

```
In [ ]: trip_records = Merged.copy()
```

```
In [ ]: trip_records["route_type"] = trip_records["route_type"].apply(lambda x:x[0])
route_to_merge = data.groupby("trip_uuid")["route_schedule_uuid"].unique().reset_index()
trip_records = trip_records.merge(route_to_merge, on="trip_uuid", how="outer")
trip_records["route_schedule_uuid"] = trip_records["route_schedule_uuid"].apply(lambda x:x[0])
trip_records
```

Out[ ]:

	trip_uuid	route_type	source_city	destination_city	source_city_state	destination_city_state
0	trip-153671041653548748	FTL	[Bhopal, Kanpur]	[Kanpur, Gurgaon]	[Bhopal Madhya Pradesh, Kanpur Uttar Pradesh]	[Kanpur Gur...
1	trip-153671042288605164	Carting	[Tumkur, Doddablpur]	[Doddablpur, Chikblapur]	[Tumkur Karnataka, Doddablpur Karnataka]	Karnata...
2	trip-153671043369099517	FTL	[Bengaluru, Gurgaon]	[Gurgaon, Chandigarh]	[Bengaluru Karnataka, Gurgaon Haryana]	[Gur...
3	trip-153671046011330457	Carting	[Mumbai]	[Mumbai]	[Mumbai Hub Maharashtra]	[Mumbai...
4	trip-153671052974046625	FTL	[Bellary, Hospet, Sandur]	[Hospet, Sandur, Bellary]	[Bellary Karnataka, Hospet Karnataka, Sandur K...	[Hos...
...	...	...	...	...	...	...
14812	trip-153861095625827784	Carting	[Chandigarh]	[Zirakpur, Chandigarh]	[Chandigarh Punjab, Chandigarh Chandigarh]	[Zir...
14813	trip-153861104386292051	Carting	[FBD]	[Faridabad]	[FBD Haryana]	[Farida...
14814	trip-153861106442901555	Carting	[Kanpur]	[Kanpur]	[Kanpur Uttar Pradesh]	[Kanpur I...
14815	trip-153861115439069069	Carting	[Tirunelveli, Eral, Tirchchndr, Thisayanvilai, Thisayanvilai,...]	[Eral, Tirchchndr, Thisayanvilai, Peikulam, Ti...]	[Tirunelveli Tamil Nadu, Eral Tamil Nadu, Tirc...	[Er...
14816	trip-153861118270144424	FTL	[Hospet, Sandur]	[Sandur, Bellary]	[Hospet Karnataka, Sandur Karnataka]	[Sand...

14817 rows × 18 columns

In [ ]: trip\_records.isna().sum()

```
Out[ ]: trip_uuid          0
         route_type        0
         source_city        0
         destination_city   0
         source_city_state  0
         destination_city_state 0
         source_state        0
         destination_state   0
         segment_osrm_time  0
         osrm_time           0
         segment_actual_time 0
         actual_time          0
         time_taken_btwn_odstart_and_od_end 0
         start_scan_to_end_scan 0
         segment_osrm_distance 0
         actual_distance_to_destination 0
         osrm_distance        0
         route_schedule_uuid  0
dtype: int64
```

## Unnesting Data

```
In [ ]: trip_records["source_city"] = trip_records["source_city"].astype("str").str.strip()
trip_records["destination_city"] = trip_records["destination_city"].astype("str").str.strip()
trip_records["source_city_state"] = trip_records["source_city_state"].astype("str")
trip_records["destination_city_state"] = trip_records["destination_city_state"].astype("str")

trip_records["source_state"] = trip_records["source_state"].astype("str").str.strip()
trip_records["destination_state"] = trip_records["destination_state"].astype("str")
```

```
In [ ]: trip_records.corr()
```

```
<ipython-input-724-5fac24159788>:1: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future version, it will default to False. Select only valid columns or specify the value of numeric_only to silence this warning.
```

```
    trip_records.corr()
```

```
Out[ ]:
```

	segment_osrm_time	osrm_time	segment_actual_time	actual_time
<b>segment_osrm_time</b>	1.000000	0.993508	0.953039	0.953039
<b>osrm_time</b>	0.993508	1.000000	0.957747	0.957747
<b>segment_actual_time</b>	0.953039	0.957747	1.000000	0.999920
<b>actual_time</b>	0.953800	0.958613	0.999920	1.000000
<b>time_taken_btwn_odstart_and_od_end</b>	0.918447	0.926280	0.961096	0.961096
<b>start_scan_to_end_scan</b>	0.918493	0.926469	0.961107	0.961107
<b>segment_osrm_distance</b>	0.996092	0.991848	0.956106	0.956106
<b>actual_distance_to_destination</b>	0.987627	0.993556	0.953048	0.953048
<b>osrm_distance</b>	0.992050	0.997610	0.958341	0.958341

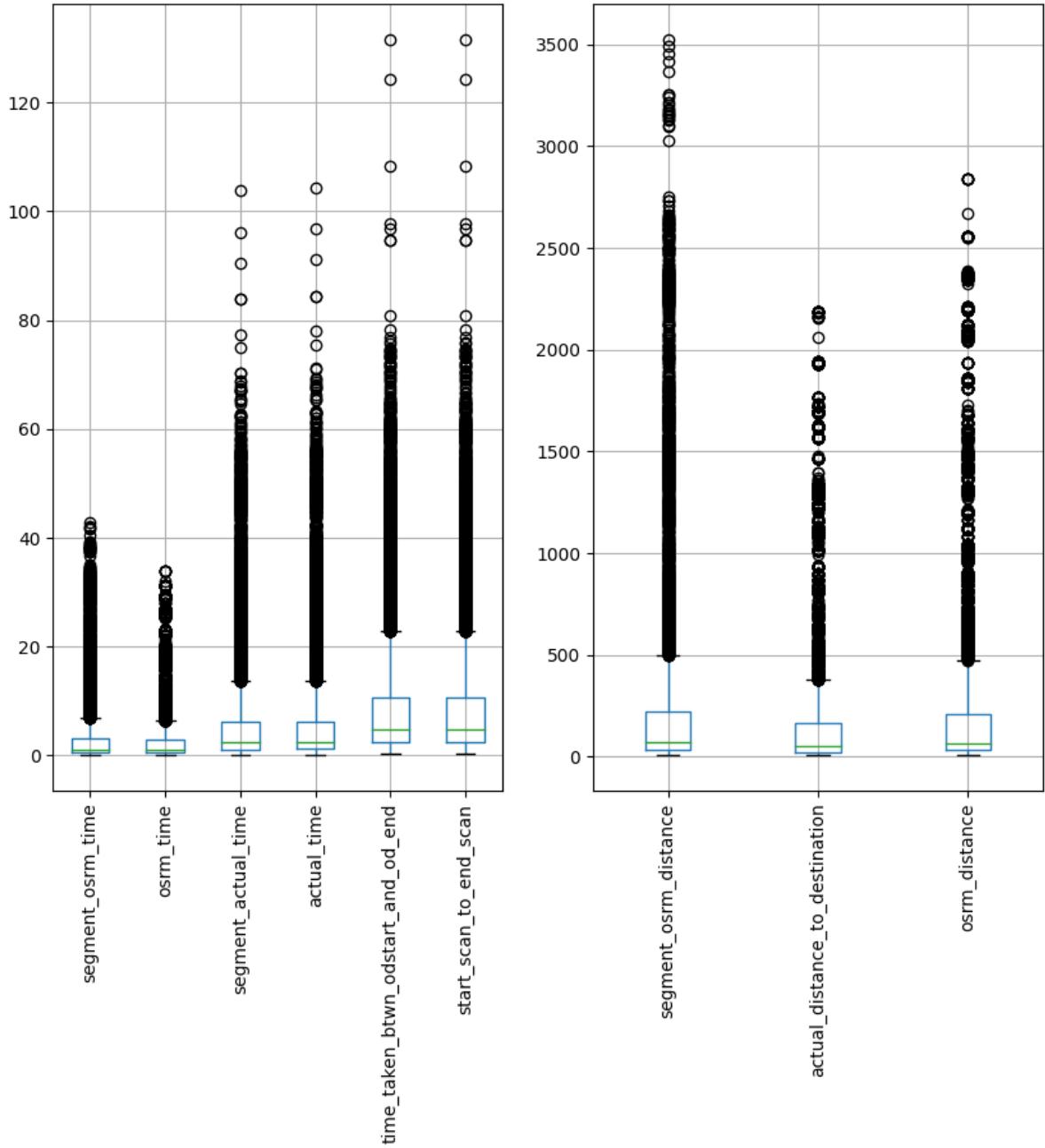
## Detecting Outliers

```
In [ ]: plt.figure(figsize = (10,8))
plt.subplot(121)
trip_records[['segment_osrm_time', 'osrm_time',
```

```

'segment_actual_time', 'actual_time',
'time_taken_btwn_odstart_and_od_end', 'start_scan_to_end_scan']].boxplot()
plt.xticks(rotation =90)
plt.subplot(122)
trip_records[['segment_osrm_distance', 'actual_distance_to_destination',
              'osrm_distance']].boxplot()
plt.xticks(rotation =90)
plt.show()

```



```
In [ ]: outlier_treatment = trip_records.copy()
```

```
In [ ]: outlier_treatment_num = outlier_treatment[['segment_osrm_time', 'osrm_time',
                                                 'segment_actual_time', 'actual_time',
                                                 'time_taken_btwn_odstart_and_od_end', 'start_scan_to_end_scan',
                                                 'segment_osrm_distance', 'actual_distance_to_destination',
                                                 'osrm_distance']]
```

## Treating Outliers

```
In [ ]: trip_records_without_outliers = trip_records.loc[outlier_treatment_num[(np.abs(stats
```

Out[ ]:

	trip_uuid	route_type	source_city	destination_city	source_city_state	destination_city_state
0	trip-153671041653548748	FTL	Bhopal Kanpur	Kanpur Gurgaon	Bhopal Madhya Pradesh Kanpur Uttar Pradesh	Kanpur U Gurgaon
1	trip-153671042288605164	Carting	Tumkur Doddablpur	Doddablpur Chikblapur	Tumkur Karnataka Doddablpur Karnataka	Doddablpur Chikblapur
3	trip-153671046011330457	Carting	Mumbai	Mumbai	Mumbai Hub Maharashtra	Mumbai Hub
4	trip-153671052974046625	FTL	Bellary Hospet Sandur	Hospet Sandur Bellary	Bellary Karnataka Hospet Karnataka Sandur Karnataka	Hospet Sandur B
5	trip-153671055416136166	Carting	Chennai	Chennai	Chennai Tamil Nadu	Chennai
...						
14812	trip-153861095625827784	Carting	Chandigarh	Zirakpur Chandigarh	Chandigarh Punjab Chandigarh Chandigarh	Zira Chandigarh
14813	trip-153861104386292051	Carting	FBD	Faridabad	FBD Haryana	Faridabad
14814	trip-153861106442901555	Carting	Kanpur	Kanpur	Kanpur Uttar Pradesh	Kanpur U
14815	trip-153861115439069069	Carting	Tirunelveli Eral Tirchchndr Thisayanvilai Peikulam Tirunel...	Eral Tirchchndr Thisayanvilai Peikulam Tirunel...	Tirunelveli Tamil Nadu Eral Tamil Nadu Tirchch...	Eral Tirchchndr
14816	trip-153861118270144424	FTL	Hospet Sandur	Sandur Bellary	Hospet Karnataka Sandur Karnataka	Sandur Bellary

14160 rows × 18 columns

### Processing Data for One hot encoding :

merging locations details into one columns and re categorise the data as per highest trips having location as top category

```
In [ ]: trip_records_without_outliers["destination_source_locations"] = trip_records_without_outliers.drop(["source_city_state","destination_city_state"],axis=1)
```

```
In [ ]: sc_dc = trip_records_without_outliers.groupby(["destination_source_locations"])["trips"].sum()
```

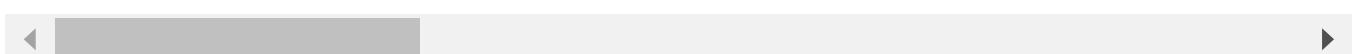
```
In [ ]: def get_cat(H):
    if 0 <= H <= 50:
        return "Category 7"
    elif 51 <= H <= 100:
        return "Category 6"
    elif 101 <= H <= 200:
        return "Category 5"
    elif 201 <= H <= 300:
        return "Category 4"
    elif 301 <= H <= 400:
        return "Category 3"
    elif 401 <= H <= 500:
        return "Category 2"
    else:
        return "Category 1"
```

```
In [ ]: sc_dc["city"] = pd.Series(map(get_cat,sc_dc["trip_uuid"]))
trip_records_for_encoding = sc_dc.merge(trip_records_without_outliers,
                                         on="destination_source_locations")
trip_records_for_encoding.drop(["destination_source_locations","trip_uuid_x"],axis=1)
trip_records_for_encoding.drop(["trip_uuid_y"],axis = 1,inplace=True)
# trip_records_for_encoding.sample(15)
encoded_data = pd.get_dummies(trip_records_for_encoding,
                               columns=["route_type","city"])
encoded_data
```

Out[ ]:

	source_city	destination_city	source_state	destination_state	segment_osrm_time	osrm_tin
0	Bengaluru	Bengaluru	Karnataka	Karnataka	1.383333	0.950000
1	Bengaluru	Bengaluru	Karnataka	Karnataka	1.150000	0.883333
2	Bengaluru	Bengaluru	Karnataka	Karnataka	1.183333	0.966667
3	Bengaluru	Bengaluru	Karnataka	Karnataka	0.700000	0.733333
4	Bengaluru	Bengaluru	Karnataka	Karnataka	0.783333	0.666667
...	...	...	...	...	...	...
14155	Hyderabad Kadthal Kalwakurthy Devarakonda	Kadthal Kalwakurthy Devarakonda Haliya	Telangana	Telangana	1.966667	1.983333
14156	Hyderabad Kadthal	Kadthal Devarakonda	Telangana	Telangana	1.483333	1.433333
14157	Hyderabad Kadthal Haliya	Kadthal Kalwakurthy Hyderabad	Telangana	Telangana	2.916667	2.866667
14158	Hyderabad Kadthal Haliya	Kadthal Devarakonda Hyderabad	Telangana	Telangana	3.383333	3.333333
14159	nan	nan	nan	nan	0.800000	0.816667

14160 rows × 23 columns



## Column Standardization

In [ ]:

```
['segment_osrm_time', 'osrm_time',
 'segment_actual_time', 'actual_time',
 'time_taken_btwn_odstart_and_od_end', 'start_scan_to_end_scan', 'segment_osr
```

Out[ ]:

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

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

```
In [ ]: scaler = StandardScaler()  
std_data = scaler.fit_transform(encoded_data[['segment_osrm_time',  
    'osrm_time',  
    'segment_actual_time',  
    'actual_time',  
    'time_taken_btwn_odstart_and_od_end',  
    'start_scan_to_end_scan',  
    'segment_osrm_distance',  
    'actual_distance_to_destination',  
    'osrm_distance']])  
std_data = pd.DataFrame(std_data, columns=['segment_osrm_time',  
    'osrm_time',  
    'segment_actual_time',  
    'actual_time',  
    'time_taken_btwn_odstart_and_od_end',  
    'start_scan_to_end_scan',  
    'segment_osrm_distance',  
    'actual_distance_to_destination',  
    'osrm_distance'])  
std_data.head()
```

```
Out[ ]:
```

	segment_osrm_time	osrm_time	segment_actual_time	actual_time	time_taken_btwn_odstart_and_od_end	start_scan_to_end_scan	segment_osrm_distance	actual_distance_to_destination	osrm_distance
0	-0.269133	-0.409683		-0.220225	-0.214843				-0.000100
1	-0.359785	-0.438916		-0.324535	-0.321822				-0.000100
2	-0.346835	-0.402374		-0.193306	-0.194785				-0.000100
3	-0.534615	-0.504692		-0.597087	-0.599297				-0.000100
4	-0.502239	-0.533926		-0.509601	-0.509034				-0.000100

```
In [ ]: scaler = MinMaxScaler()  
MinMax_data = scaler.fit_transform(encoded_data[['segment_osrm_time','osrm_time','segment_actual_time','actual_time','time_taken_btwn_odstart_and_od_end','start_scan_to_end_scan','segment_osrm_distance','osrm_distance']])  
MinMax_data = pd.DataFrame(MinMax_data,columns=['segment_osrm_time','osrm_time','segment_actual_time','actual_time','time_taken_btwn_odstart_and_od_end','start_scan_to_end_scan','segment_osrm_distance','actual_distance_to_destination','osrm_distance'])  
MinMax_data.head()
```

```
Out[ ]:
```

	segment_osrm_time	osrm_time	segment_actual_time	actual_time	time_taken_btwn_odstart_and_od_end	start_scan_to_end_scan	segment_osrm_distance	actual_distance_to_destination	osrm_distance
0	0.069369	0.059302		0.098113	0.098719				0.000100
1	0.056757	0.054651		0.081402	0.081644				0.000100
2	0.058559	0.060465		0.102426	0.101921				0.000100
3	0.032432	0.044186		0.037736	0.037353				0.000100
4	0.036937	0.039535		0.051752	0.051761				0.000100

```
In [ ]: std_data
```

Out[ ]:

	segment_osrm_time	osrm_time	segment_actual_time	actual_time	time_taken_btwn_odstart
0	-0.269133	-0.409683	-0.220225	-0.214843	
1	-0.359785	-0.438916	-0.324535	-0.321822	
2	-0.346835	-0.402374	-0.193306	-0.194785	
3	-0.534615	-0.504692	-0.597087	-0.599297	
4	-0.502239	-0.533926	-0.509601	-0.509034	
...	...	...	...	...	
14155	-0.042502	0.043440	-0.210131	-0.211500	
14156	-0.230282	-0.197738	-0.314441	-0.311792	
14157	0.326583	0.430787	0.136448	0.136179	
14158	0.507888	0.635424	1.347789	1.336342	
14159	-0.495764	-0.468150	-0.435575	-0.435486	

14160 rows × 9 columns

In [ ]:

```
one_hot_encoded_data = encoded_data[["route_type_Carting","route_type_FTL","city_Cat
"city_Category 2","city_Category 3","city_Category 4",
"city_Category 5","city_Category 6","city_Category 7"]]
```

In [ ]:

```
Standardized_Data = pd.concat([std_data,one_hot_encoded_data],axis = 1)
```

In [ ]:

```
Min_Max_Scaled_Data = pd.concat([MinMax_data,one_hot_encoded_data],axis = 1)
```

In [ ]:

```
Standardized_Data.sample(5)
```

Out[ ]:

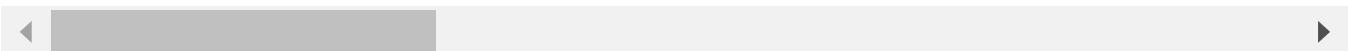
	segment_osrm_time	osrm_time	segment_actual_time	actual_time	time_taken_btwn_odstart
10715	-0.197906	-0.131962	-0.011605	-0.017602	
12569	0.300683	0.386937	0.133083	0.132836	
11101	-0.683544	-0.680095	-0.741775	-0.739706	
11055	-0.599366	-0.577777	-0.674478	-0.672844	
567	-0.612317	-0.614319	-0.741775	-0.743049	

In [ ]:

```
Min_Max_Scaled_Data.sample(5)
```

Out[ ]:

	segment_osrm_time	osrm_time	segment_actual_time	actual_time	time_taken_btwn_odstart
1727	0.022523	0.029070	0.024259	0.024546	
30	0.015315	0.019767	0.014555	0.014408	
12977	0.166667	0.175581	0.171968	0.170224	
1174	0.018919	0.026744	0.016173	0.016542	
5494	0.025225	0.032558	0.031267	0.031483	



Route analysis :

```
In [ ]: A = data.groupby("route_schedule_uuid")["route_type"].unique().reset_index()
B = data.groupby("route_schedule_uuid")["destination_city"].unique().reset_index()
B.columns = ["route_schedule_uuid","destination_cities"]
C = data.groupby("route_schedule_uuid")["source_city"].unique().reset_index()
C.columns = ["route_schedule_uuid","source_cities"]
D = data.groupby("route_schedule_uuid")["source_state"].unique().reset_index()
D.columns = ["route_schedule_uuid","source_states"]
E = data.groupby("route_schedule_uuid")["destination_state"].unique().reset_index()
E.columns = ["route_schedule_uuid","destination_states"]
F = data.groupby("route_schedule_uuid")[[ "source_state",
                                         "destination_state"]].nunique().sort_values
                                         F.columns = ["route_schedule_uuid","#source_states"
                                         , "#destination_states"]
G = trip_records.groupby("route_schedule_uuid")["actual_distance_to_destination"].nunique()
G.columns = ["route_schedule_uuid","Average_Actual_distance_to_destination"]
H = trip_records["route_schedule_uuid"].value_counts().reset_index()
H.columns = ["route_schedule_uuid","Number_of_Trips"]
I = data.groupby("route_schedule_uuid")[[ "source_city",
                                         "destination_city"]].nunique().sort_values
                                         I.columns = ["route_schedule_uuid","#source_cities"
                                         , "#destination_cities"]
```

```
In [ ]: route_records = I.merge(H.merge(G.merge(F.merge(E.merge(D.merge(C.merge(A.merge(B,
                                         on ="route_schedule_uuid",
                                         how = "outer"),on ="route_schedule_uuid",
                                         how = "outer"),
                                         on ="route_schedule_uuid",
                                         how = "outer")
```

```
In [ ]: route_records.isna().sum()
```

```
Out[ ]: route_schedule_uuid          0
         #source_cities            0
         #destination_cities       0
         Number_of_Trips          0
         Average_Actual_distance_to_destination 0
         #source_states           0
         #destination_states      0
         destination_states       0
         source_states             0
         source_cities              0
         route_type                 0
         destination_cities        0
dtype: int64
```

```
In [ ]: route_records.dropna(inplace=True)
```

```
In [ ]: route_records["route_type"] = route_records["route_type"].astype("str").str.strip(' ')
route_records["source_cities"] = route_records["source_cities"].astype("str").str.strip(' ')
route_records["destination_cities"] = route_records["destination_cities"].astype("str").str.strip(' ')
route_records["source_states"] = route_records["source_states"].astype("str").str.strip(' ')
route_records["destination_states"] = route_records["destination_states"].astype("str").str.strip(' ')
```

```
In [ ]: route_records
```

Out[ ]:

	route_schedule_uuid	#source_cities	#destination_cities	Number_of_Trips	Average_Actual_
--	---------------------	----------------	---------------------	-----------------	-----------------

<b>0</b>	thanos::sroute:d010efca-d90d-4977-b987-eae68c5...	13	11	14
<b>1</b>	thanos::sroute:4cbecb35-356b-4b68-bf3c-6225b5e...	10	10	12
<b>2</b>	thanos::sroute:ae5c430f-6153-48d1-8fe5-d5f0bbc...	10	10	20
<b>3</b>	thanos::sroute:f8968c72-5222-4d81-9eed-8a6d88f...	9	9	9
<b>4</b>	thanos::sroute:ed5b80be-7abf-424d-b8cd-d81556a...	9	8	20
...	...	...	...	...
<b>1499</b>	thanos::sroute:9e7bb811-593f-47bc-ac49-ba03ed8...	1	1	19
<b>1500</b>	thanos::sroute:46b9641b-55b5-4b15-b039-2612a50...	1	1	15
<b>1501</b>	thanos::sroute:b48f633d-15cb-4744-a0b9-21df0a9...	1	1	7
<b>1502</b>	thanos::sroute:265efe06-3625-4fba-afee-07b5b64...	0	1	1
<b>1503</b>	thanos::sroute:cfd575b8-df26-48f5-8427-6f48f9d...	0	0	1

1504 rows × 12 columns



In [ ]:

```
route_records["ROUTE"] = route_records["source_cities"] + " -- " + route_records["destination_cities"]
route_records.drop(["route_schedule_uuid"], axis = 1, inplace=True)
first_column = route_records.pop('ROUTE')
route_records.insert(0, 'ROUTE', first_column)
route_records["SourceToDestination_city"] = route_records["source_cities"].str.split(" -- ")
```

```

first_column = route_records.pop('SouceToDestination_city')
route_records.insert(0, 'SouceToDestination_city', first_column)
route_records

```

Out[ ]:

	SouceToDestination_city	ROUTE	#source_cities	#destination_cities	Number_of_Trips	Average_Distance
0	Guwahati TO LakhimpurN	Guwahati LakhimpurN Dhemaji Likabali Tezpur Pa...	13	11	11	14
1	Guwahati TO Tura	Guwahati Rangia Kokrajhar Dhubri Bilasipara Tu...	10	10	10	12
2	Jaipur TO Tarnau	Jaipur Chomu Reengus Sikar Bikaner Didwana Suj...	10	10	10	20
3	Mangalore TO Udupi	Mangalore Udupi Kundapura Bhatkal Honnavar Kum...	9	9	9	9
4	Ajmer TO Raipur	Ajmer Beawar Bilara Bijainagar Kekri Nasirabad...	9	8	8	20
...	...	...	...	...	...	...
1499	Mumbai TO Mumbai	Mumbai -- Mumbai	1	1	1	19
1500	Mumbai TO Mumbai	Mumbai -- Mumbai	1	1	1	15
1501	Bengaluru TO Bengaluru	Bengaluru - - Bengaluru	1	1	1	7
1502	nan TO Mainpuri	nan -- Mainpuri	0	1	1	1
1503	nan TO nan	nan -- nan	0	0	0	1

1504 rows × 13 columns

Exploratory Data Analysis : ( getting some insights from preprocessed data ) :

Busiest Route Analysis : Number of Trips between cities , sorted highest to lowest

Top 20 source and destination cities which have high frequency of trips in between .

```
In [ ]: Number_of_trips_between_cities = data.groupby(["source_city_state", "destination_city_state"])["trip_uuid"].count()  
Number_of_trips_between_cities.head(25)
```

```
Out[ ]:
```

	source_city_state	destination_city_state	trip_uuid
0	Bengaluru Karnataka	Bengaluru Karnataka	1369
1	Bhiwandi Maharashtra	Mumbai Maharashtra	512
2	Mumbai Maharashtra	Mumbai Maharashtra	361
3	Hyderabad Telangana	Hyderabad Telangana	308
4	Mumbai Maharashtra	Bhiwandi Maharashtra	282
5	Delhi Delhi	Gurgaon Haryana	248
6	Gurgaon Haryana	Delhi Delhi	237
7	Mumbai Hub Maharashtra	Mumbai Maharashtra	227
8	Chennai Tamil Nadu	Chennai Tamil Nadu	205
9	MAA Tamil Nadu	Chennai Tamil Nadu	204
10	Chennai Tamil Nadu	MAA Tamil Nadu	141
11	Bengaluru Karnataka	HBR Karnataka	133
12	Ahmedabad Gujarat	Ahmedabad Gujarat	131
13	Pune Maharashtra	PNQ Maharashtra	122
14	Jaipur Rajasthan	Jaipur Rajasthan	111
15	Delhi Delhi	Delhi Delhi	109
16	Pune Maharashtra	Bhiwandi Maharashtra	107
17	Pune Maharashtra	Pune Maharashtra	101
18	Chandigarh Chandigarh	Chandigarh Punjab	100
19	Kolkata West Bengal	CCU West Bengal	96
20	Gurgaon Haryana	Sonipat Haryana	92
21	Sonipat Haryana	Gurgaon Haryana	86
22	Chandigarh Punjab	Chandigarh Chandigarh	84
23	HBR Karnataka	Bengaluru Karnataka	79
24	Bengaluru Karnataka	BLR Karnataka	78

From above table, we can observe that Mumbai Maharashtra ,Delhi ,Gurgaon(Haryana),Bengaluru Karnataka ,Hyderabad Telangana,Chennai Tamil Nadu,Ahmedabad Gujarat,Pune Maharashtra,Chandigarh Chandigarh and Kolkata West Bengal are some cities have higest amount of trips happening states with in the city.

```
In [ ]: Number_of_trips_between_cities.loc[Number_of_trips_between_cities["source_city_stat"]]
```

Out[ ]:

	source_city_state	destination_city_state	trip_uuid
1	Bhiwandi Maharashtra	Mumbai Maharashtra	512
4	Mumbai Maharashtra	Bhiwandi Maharashtra	282
5	Delhi Delhi	Gurgaon Haryana	248
6	Gurgaon Haryana	Delhi Delhi	237
7	Mumbai Hub Maharashtra	Mumbai Maharashtra	227
9	MAA Tamil Nadu	Chennai Tamil Nadu	204
10	Chennai Tamil Nadu	MAA Tamil Nadu	141
11	Bengaluru Karnataka	HBR Karnataka	133
13	Pune Maharashtra	PNQ Maharashtra	122
16	Pune Maharashtra	Bhiwandi Maharashtra	107
18	Chandigarh Chandigarh	Chandigarh Punjab	100
19	Kolkata West Bengal	CCU West Bengal	96
20	Gurgaon Haryana	Sonipat Haryana	92
21	Sonipat Haryana	Gurgaon Haryana	86
22	Chandigarh Punjab	Chandigarh Chandigarh	84
23	HBR Karnataka	Bengaluru Karnataka	79
24	Bengaluru Karnataka	BLR Karnataka	78
26	Del Delhi	Gurgaon Haryana	76
27	Bhiwandi Maharashtra	Pune Maharashtra	72
28	Ludhiana Punjab	Chandigarh Punjab	71
30	Chandigarh Punjab	Gurgaon Haryana	66
31	Gurgaon Haryana	Bengaluru Karnataka	66
32	LowerParel Maharashtra	Mumbai Maharashtra	65
34	Mumbai Hub Maharashtra	Bhiwandi Maharashtra	63
35	PNQ Maharashtra	Pune Maharashtra	62

For trips between different states and cities:

The highest number of trips occur between Delhi and Gurgaon. There is significant traffic between Gurgaon, Haryana, and Bengaluru, Karnataka. Bhiwandi/Mumbai, Maharashtra, to Pune, Maharashtra, is another common route. Trips from Sonipat to Gurgaon, Haryana, are also frequent. Deliveries to airports:

There is a notable volume of deliveries to airports, such as: Chennai to Chennai International Airport (MAA) Pune to Pune Airport (PNQ) Kolkata to Kolkata International Airport (CCU) in West Bengal Bengaluru to Bengaluru International Airport (BLR)

In [ ]:

```
route_records[[ "ROUTE", "Number_of_Trips",
                 "Average_Actual_distance_to_destination",
```

```
#source_cities",
"#destination_cities"]].sort_values(by="Number_of_Trips",ascending=False)
```

Out[ ]:

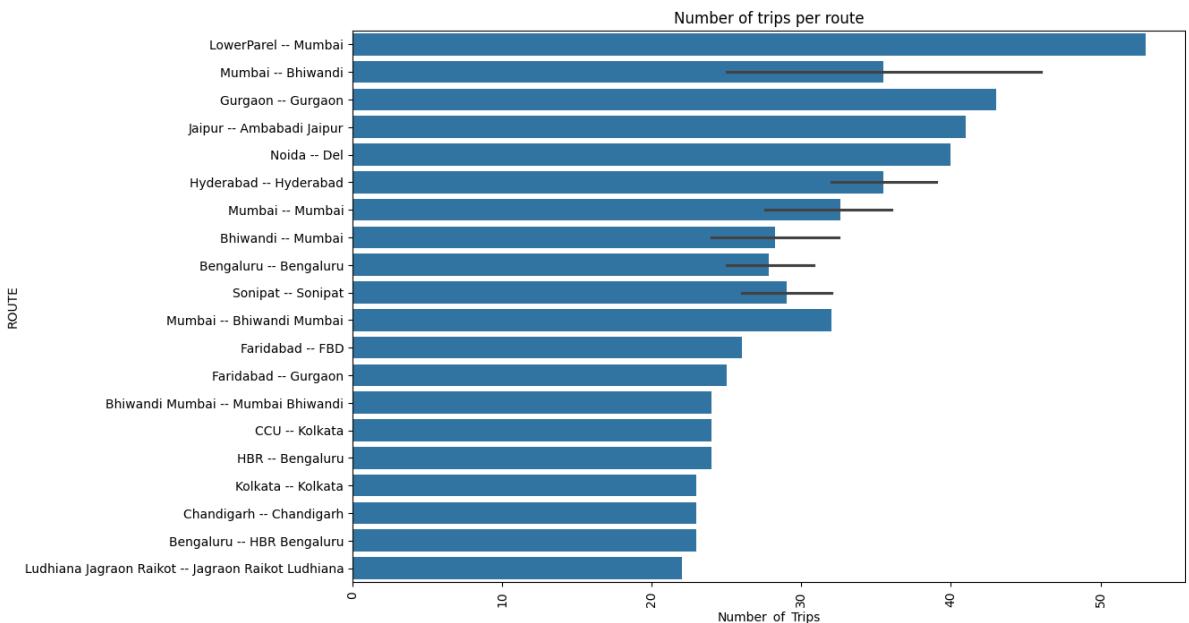
	ROUTE	Number_of_Trips	Average_Actual_distance_to_destination	#source_cities	#destina
1465	LowerParel -- Mumbai	53	16.428868	1	
1426	Mumbai -- Bhiwandi	46	20.199445	1	
808	Gurgaon -- Gurgaon	43	29.740842	1	
679	Jaipur -- Ambabadi Jaipur	41	15.348495	1	
1257	Noida -- Del	40	10.882902	1	
1368	Hyderabad -- Hyderabad	39	35.695641	1	
1273	Mumbai -- Mumbai	37	13.882863	1	
1359	Mumbai -- Mumbai	36	17.526251	1	
1303	Bhiwandi -- Mumbai	35	21.241534	1	
700	Mumbai -- Mumbai	34	15.906614	1	
751	Mumbai -- Mumbai	33	15.668726	1	
1060	Bengaluru -- Bengaluru	33	28.067004	1	
793	Sonipat -- Sonipat	32	11.691243	1	
972	Hyderabad -- Hyderabad	32	21.835579	1	
1184	Mumbai -- Bhiwandi Mumbai	32	21.601109	1	
874	Bengaluru -- Bengaluru	30	28.055789	1	
1177	Bhiwandi -- Mumbai	30	21.396002	1	
1354	Bengaluru -- Bengaluru	27	27.967087	1	
921	Faridabad -- FBD	26	9.677121	1	
1480	Sonipat -- Sonipat	26	12.182486	1	

	ROUTE	Number_of_Trips	Average_Actual_distance_to_destination	#source_cities	#destina
1041	Mumbai -- Bhiwandi	25	19.942191	1	
877	Faridabad -- Gurgaon	25	47.091622	1	
833	Bhiwandi -- Mumbai	25	21.531705	1	
1249	Bengaluru -- Bengaluru	25	28.019668	1	
869	Bengaluru -- Bengaluru	24	41.396497	1	

Top Routes having Maximum Number of Trips between/within the source and destinations

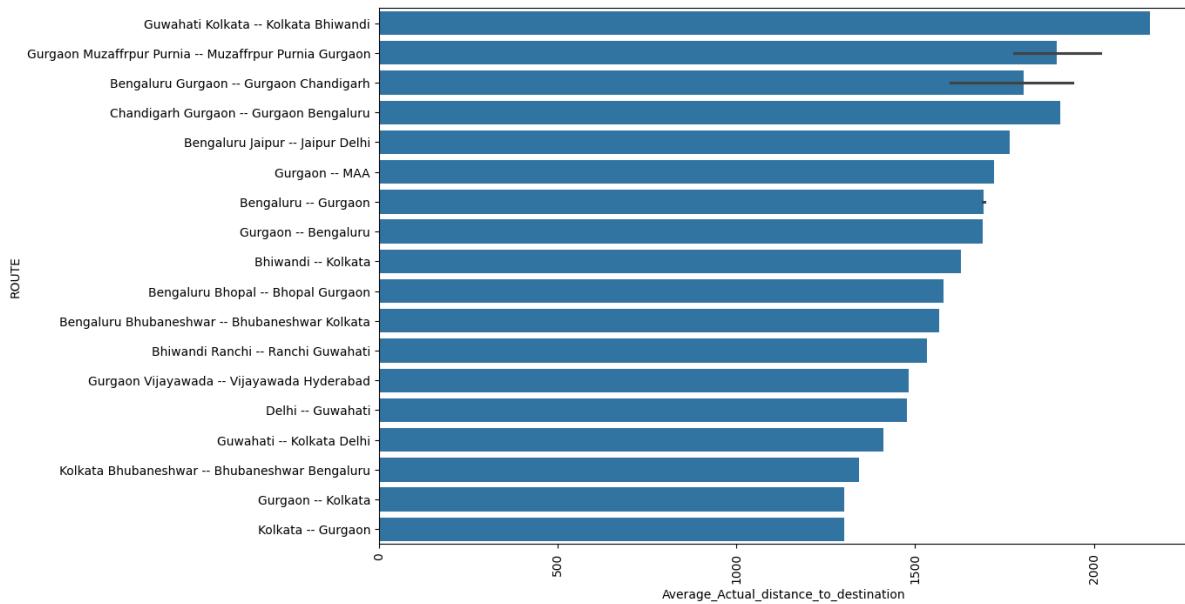
```
In [ ]: plt.figure(figsize=(12,8))

X = route_records[['ROUTE', 'Number_of_Trips',
                   ]].sort_values(by="Number_of_Trips", ascending=False).head(35)
sns.barplot(y = X["ROUTE"],
            x= X["Number_of_Trips"])
plt.title("Number of trips per route")
plt.xticks(rotation = 90)
plt.show()
```



```
In [ ]: plt.figure(figsize=(12,8))

X = route_records[['ROUTE', 'Average_Actual_distance_to_destination',
                   ]].sort_values(by="Average_Actual_distance_to_destination", ascending=True)
sns.barplot(y = X["ROUTE"],
            x = X["Average_Actual_distance_to_destination"])
plt.xticks(rotation = 90)
plt.show()
```



From above Bar chart, and table , we can observe that highest trips are happening is with in the particular cities.

In terms of average distance between destinations , we can observe Guwahati to Mumbai , Bangalore to Chandigarh ,Bangalore to Delhi , Bangalore to Gurgaon are the longest routes .

### Busiest and Longest Routes :

```
In [ ]: Busiest_and_Longest_Routes = route_records[(route_records["Average_Actual_distance"]
                                             & (route_records["Number_of_Trips"] > route_records["Number_of_Trips"]))

Busiest_and_Longest_Routes_top25 = Busiest_and_Longest_Routes[[ "source_cities",
                                                               "destination_cities",
                                                               "Number_of_Trips",
                                                               "Average_Actual_dist
Busiest_and_Longest_Routes_top25
```

Out[ ]:

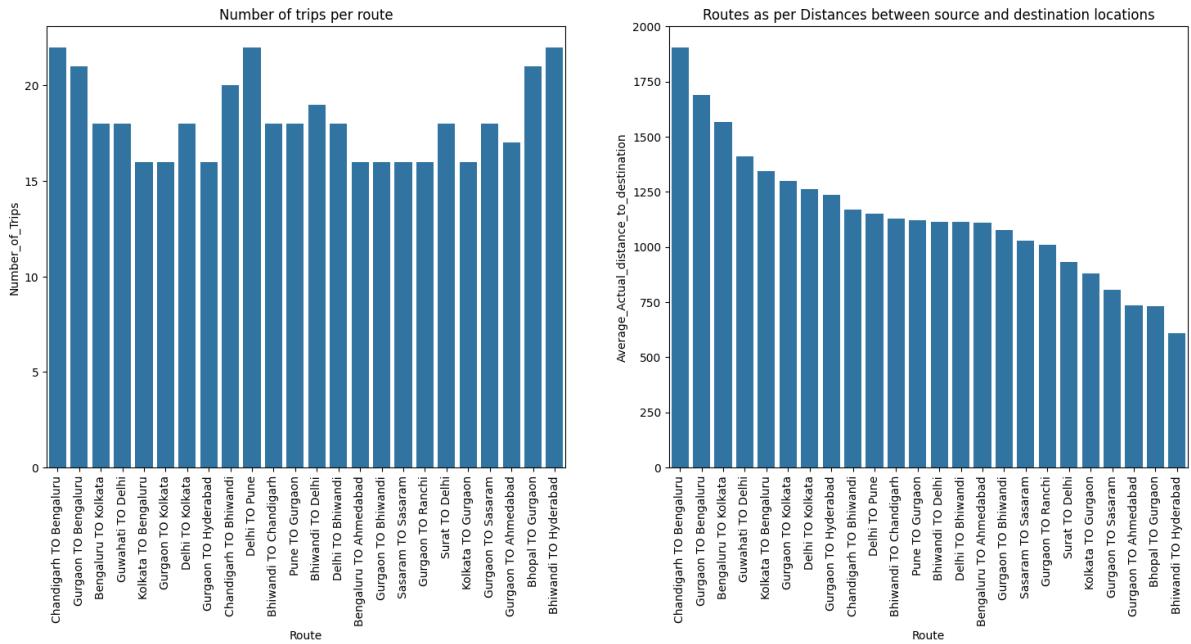
	source_cities	destination_cities	Number_of_Trips	Average_Actual_distance_to_destination
629	Chandigarh Gurgaon	Gurgaon Bengaluru	22	1905.766051
995	Gurgaon	Bengaluru	21	1689.873158
991	Gurgaon	Bengaluru	21	1689.791894
512	Bengaluru Bhubaneshwar	Bhubaneshwar Kolkata	18	1567.577507
745	Guwahati	Kolkata Delhi	18	1411.208424
624	Kolkata Bhubaneshwar	Bhubaneshwar Bengaluru	16	1342.143081
752	Gurgaon	Kolkata	16	1300.572161
588	Delhi Gurgaon	Gurgaon Kolkata	18	1263.113211
826	Gurgaon	Hyderabad	16	1236.572072
541	Chandigarh Gurgaon	Gurgaon Bhiwandi	20	1170.817927
442	Delhi Gurgaon	Gurgaon Pune	22	1151.514940
445	Bhiwandi Sonipat	Sonipat Chandigarh	18	1129.609705
739	Pune	Gurgaon	18	1120.729446
1377	Bhiwandi	Delhi	19	1114.214670
1049	Delhi	Bhiwandi	18	1114.182197
313	Bengaluru Kolhapur Surat	Kolhapur Surat Ahmedabad	16	1110.015339
1219	Gurgaon	Bhiwandi	16	1078.076312
197	Sasaram Kanpur Kolkata Dhanbad	Kanpur Gurgaon Dhanbad Sasaram	16	1028.024726
1136	Gurgaon	Ranchi	16	1010.953223
1286	Surat	Delhi	18	931.980821
439	Kolkata Ranchi	Ranchi Gurgaon	16	881.621264
1108	Gurgaon	Sasaram	18	804.210670
1454	Gurgaon	Ahmedabad	17	735.550450
223	Bhopal Kanpur Auraiya Etawah	Kanpur Auraiya Etawah Gurgaon	21	731.634456
863	Bhiwandi	Hyderabad	22	607.514619

Above Table shows the source to destination city routes having largest numbers of trip happening having large distances : which are :

Chandigarh TO Bengaluru Gurgaon TO Bengaluru Bengaluru TO Kolkata Guwahati TO Delhi  
 Delhi TO Kolkata Chandigarh TO Gurgaon Gurgaon TO Hyderabad Bengaluru TO Ahmedabad  
 Surat TO Delhi Gurgaon TO Ahmedabad

```
In [ ]: Busiest_and_Longest_Routes_top25["Route"] = Busiest_and_Longest_Routes_top25["source_cities"] + " TO " + Busiest_and_Longest_Routes_top25["destination_cities"]
Busiest_and_Longest_Routes_top25.drop(["source_cities","destination_cities"],axis = 1,inplace=True)
plt.figure(figsize=(18,7))

plt.subplot(121)
plt.title("Number of trips per route")
sns.barplot(x=Busiest_and_Longest_Routes_top25["Route"],
            y = Busiest_and_Longest_Routes_top25["Number_of_Trips"])
plt.xticks(rotation = 90)
plt.subplot(122)
plt.title("Routes as per Distances between source and destination locations")
sns.barplot(x=Busiest_and_Longest_Routes_top25["Route"],
            y= Busiest_and_Longest_Routes_top25["Average_Actual_distance_to_destination"])
plt.xticks(rotation = 90)
plt.show()
```



Routes : passing through maximum number of cities :

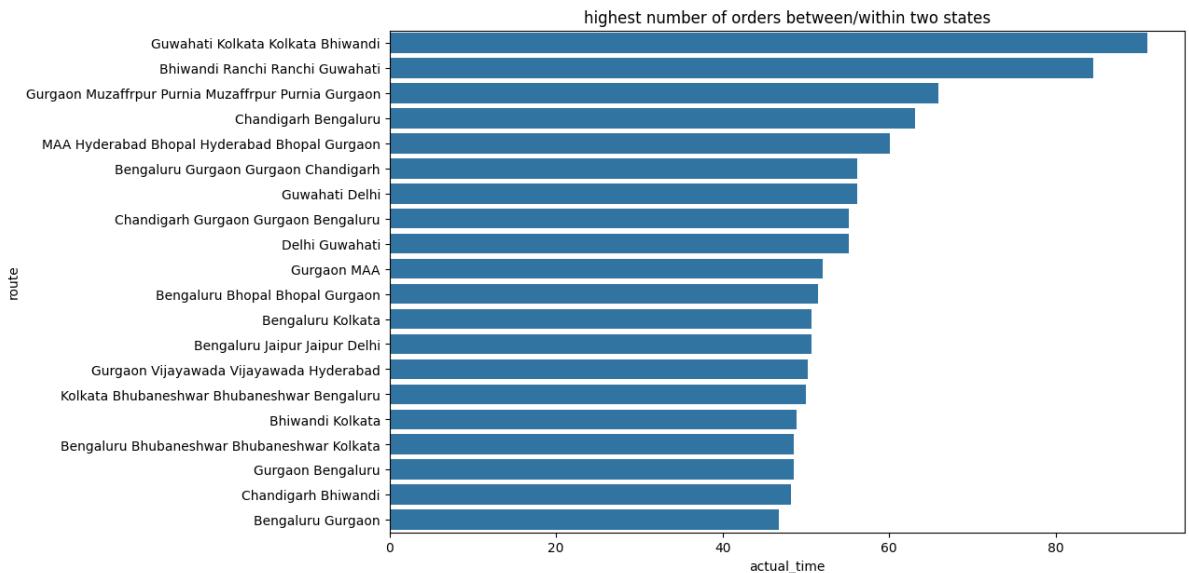
```
In [ ]: route_records[["SourceToDestination_city","Number_of_Trips",
                    "Average_Actual_distance_to_destination",
                    "#source_cities",
                    "#destination_cities"]].sort_values(by=["#source_cities",
                    "#destination_cities",
                    "Number_of_Trips"]
                    ,ascending=False).head(25)
```

	SourceToDestination_city	Number_of_Trips	Average_Actual_distance_to_destination	#source_cities
0	Guwahati TO LakhimpurN	14	281.596486	1
2	Jaipur TO Tarnau	20	351.611796	1
1	Guwahati TO Tura	12	332.602225	1
3	Mangalore TO Udupi	9	195.257193	1
4	Ajmer TO Raipur	20	178.737233	1
5	Mainpuri TO Tilhar	12	207.247057	1
8	Hassan TO Koppa	21	200.497832	1
15	Shrirampur TO Sangamner	20	204.509529	1
7	Musiri TO Tiruchi	19	219.845121	1
9	Bijnor TO Bijnor	17	209.400685	1
10	Dausa TO Lalsot	17	232.408310	1
17	Tinusukia TO Dibrugarh	16	111.098543	1
12	Pondicherry TO Pondicherry	12	230.253602	1
14	Mysore TO Mysore	12	154.324190	1
6	Golaghat TO Guwahati	11	258.546587	1
13	Varanasi TO Varanasi	8	82.545019	1
16	Vijayawada TO Suryapet	8	407.029391	1
11	Hyderabad TO Miryalguda	7	420.603709	1
27	Srikakulam TO Bobbili	22	154.495283	1
36	Pukhrayan TO Kanpur	22	139.834945	1
48	Dhule TO Shirpur	22	150.016233	1
30	Madhupur TO Madhupur	21	252.072259	1
38	Kamareddy TO Kamareddy	21	177.923330	1
42	Noida TO Khurja	21	208.714043	1
20	Junagadh TO Veraval	19	179.538596	1

Top 20 Longest Route as per : average actual time taken from one city to another city :

```
In [ ]: Longest_route_as_per_actual_trip_time = trip_records.groupby(["source_city",
                                                               "destination_city"])["actual_time"].mean().sort_values(ascending=True)
Longest_route_as_per_actual_trip_time["route"] = Longest_route_as_per_actual_trip_time.index
Longest_route_as_per_actual_trip_time.drop(["source_city",
                                             "destination_city"], axis = 1,inplace=True)
Longest_route_as_per_actual_trip_time
plt.figure(figsize=(11,7))
sns.barplot(y = Longest_route_as_per_actual_trip_time["route"],
            x = Longest_route_as_per_actual_trip_time["actual_time"],)
```

```
plt.title("highest number of orders between/within two states")
plt.show()
```



highest number of Trips happening between/within two states :

```
In [ ]: highest_order_between_states = data.groupby([ "source_state",
                                                    "destination_state"])[ "trip_uuid"].nun
HOBS = highest_order_between_states.head(15)
HOBS[ "souce-destination" ] = HOBS[ "source_state" ] + " - " + HOBS[ "destination_state" ]
HOBS.drop([ "source_state", "destination_state"], axis = 1, inplace=True)
HOBS.columns = [ "Number_of_trips_between_states", "souce-destination_state" ]

plt.figure(figsize=(11,5))
sns.barplot(y = HOBS[ "souce-destination_state" ],
            x = HOBS[ "Number_of_trips_between_states" ],)
plt.title("highest number of orders within two states")
plt.show()
```

```
<ipython-input-759-666a8f846c08>:4: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead
```

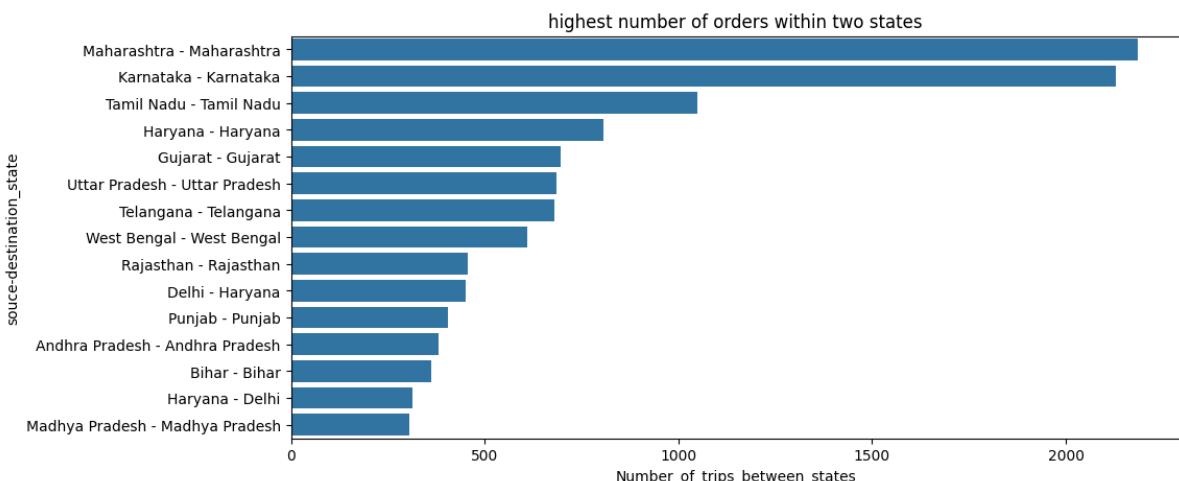
See the caveats in the documentation: [https://pandas.pydata.org/pandas-docs/stable/user\\_guide/indexing.html#returning-a-view-versus-a-copy](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)

```
HOBS[ "souce-destination" ] = HOBS[ "source_state" ] + " - " + HOBS[ "destination_state" ]
```

```
<ipython-input-759-666a8f846c08>:5: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame
```

See the caveats in the documentation: [https://pandas.pydata.org/pandas-docs/stable/user\\_guide/indexing.html#returning-a-view-versus-a-copy](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)

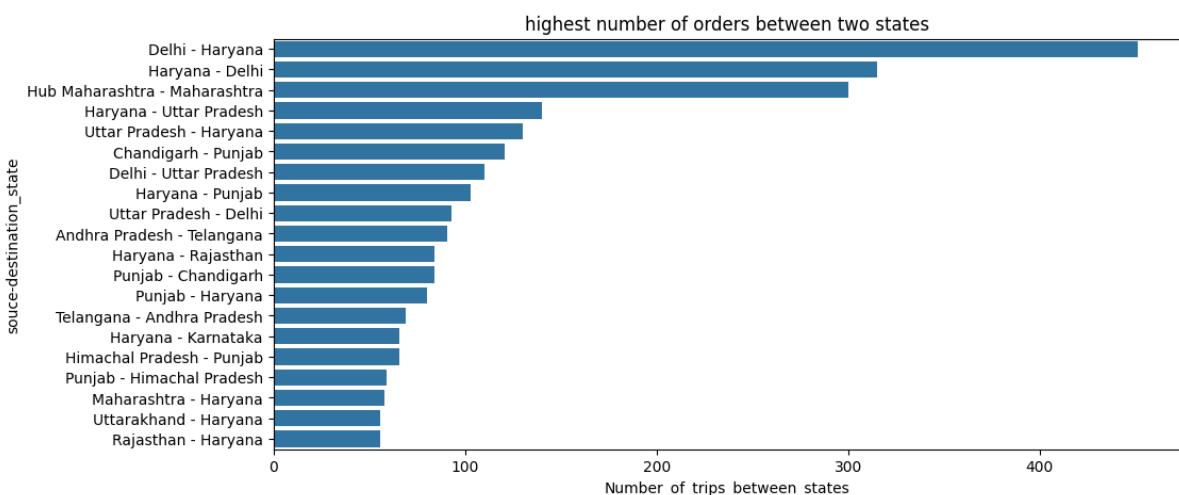
```
HOBS.drop([ "source_state", "destination_state"], axis = 1, inplace=True)
```



```
In [ ]: HOBS = data.groupby(["source_state","destination_state"])["trip_uuid"].nunique().sort_values(ascending=False)
HOBS = HOBS[HOBS["source_state"]!=HOBS["destination_state"]].head(20)

HOBS["souce-destination"] = HOBS["source_state"] + " - " + HOBS["destination_state"]
HOBS.drop(["source_state","destination_state"],axis = 1, inplace=True)
HOBS.columns = ["Number_of_trips_between_states","souce-destination_state"]

plt.figure(figsize=(11,5))
sns.barplot(y = HOBS["souce-destination_state"],
            x = HOBS["Number_of_trips_between_states"],)
plt.title("highest number of orders between two states")
plt.show()
```



From above charts ,

Delhi to Haryana is the busiest route, having more than 400 trips in between.

Some of such busy routes are Haryana to Uttar Pradesh , Chandigarh to Punjab , Delhi to Uttar Pradesh.

Within the state , Maharashtra , Karnataka, Tamil Nadu are some states having above 1000 trips.

## Top 20 warehouses with heavy traffic :

```

transactions.columns = ["source_city_state","#Trips_s","destination_city_state","#Trips_d"]
transactions["TripsTraffic"] = transactions["#Trips_s"]+transactions["#Trips_d"]
transactions.drop(["#Trips_s","#Trips_d","destination_city_state"],axis = 1,inplace=True)
transactions.columns = ["Warehouse_City(Junction)","TripsTraffic"]

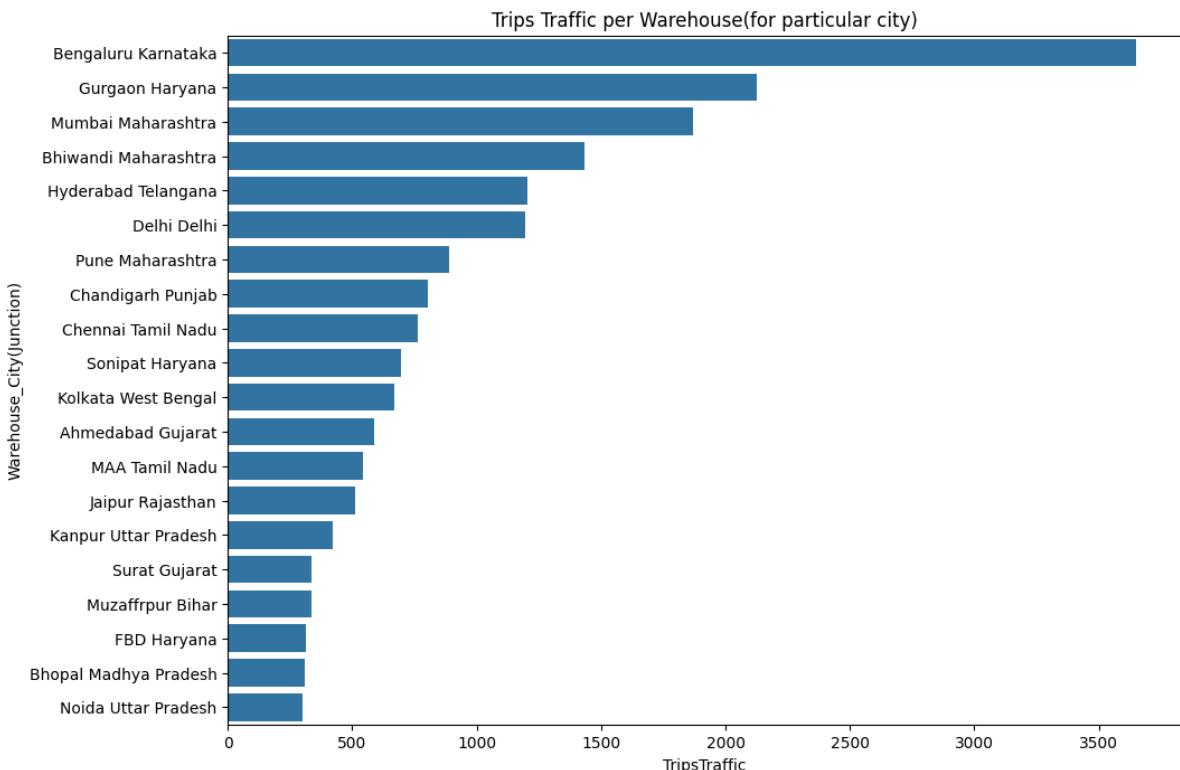
```

```
In [ ]: T = transactions.sort_values(by=["TripsTraffic"],ascending=False).head(20)
```

```

In [ ]: plt.figure(figsize=(11,8))
sns.barplot(y = T["Warehouse_City(Junction)"],
             x = T["TripsTraffic"])
plt.title("Trips Traffic per Warehouse(for particular city)")
plt.show()

```



Top 20 Busiest Warehouse (junctions) as per trips traffic at the juction : are

'Bengaluru Karnataka', 'Gurgaon Haryana', 'Mumbai Maharashtra', 'Bhiwandi Maharashtra', 'Hyderabad Telangana', 'Delhi Delhi', 'Pune Maharashtra', 'Chandigarh Punjab', - 'Chennai Tamil Nadu', 'Sonipat Haryana', - 'Kolkata West Bengal', 'Ahmedabad Gujarat', 'MAA Tamil Nadu', 'Jaipur Rajasthan', 'Kanpur Uttar Pradesh', - 'Surat Gujarat', 'Muzaffarpur Bihar', 'FBD Haryana', 'Bhopal Madhya Pradesh', 'Noida Uttar Pradesh'

```
In [ ]: trip_records.groupby(["source_state","destination_state"])["trip_uuid"].count().sort_index()
```

Out[ ]:

	source_state	destination_state	trip_uuid
0	Maharashtra	Maharashtra	2085
1	Karnataka	Karnataka	2002
2	Tamil Nadu	Tamil Nadu	996
3	Haryana	Haryana	771
4	Telangana	Telangana	627
5	Gujarat	Gujarat	624
6	West Bengal	West Bengal	610
7	Uttar Pradesh	Uttar Pradesh	529
8	Rajasthan	Rajasthan	400
9	Delhi	Haryana	385
10	Andhra Pradesh	Andhra Pradesh	344
11	Punjab	Punjab	342
12	Bihar	Bihar	330
13	Haryana	Delhi	307
14	Hub Maharashtra	Maharashtra	300

## Insights

-14817 different trips happened between source to destinations during 2018 , September and October.

-1504 delivery routes on which trips are happenig.

-we have 1508 unique source centers and 1481 unique destination centers

-From 14817 total different trips , we have 8908 (60%) of the trip-routes are Carting , which consists of small vehicles and 5909 (40%) of total trip-routes are FTL : which are Full Truck Load get to the destination sooner. as no other pickups or drop offs along the way .

## Hypothesis tests Results

- from 2 sample t-test ,we can also conclude that
- Average time\_taken\_btwn\_odstart\_and\_od\_end for population is equal to Average start\_scan\_to\_end\_scan for population.
- population average actual\_time is less than population average start\_scan\_to\_end\_scan.
- population mean Actual time taken to complete delivery and population mean time\_taken\_btwn\_od\_start\_and\_od\_end are also not same.
- Mean of actual time is higher than Mean of the OSRM estimated time for delivery
- Population average for Actual Time taken to complete delivery trip and segment actual time are same.

- Average of OSRM Time & segment-osrm-time for population is not same.
- Population Mean osrm time is less than Population Mean segment osrm time.
- Average of OSRM distance for population is less than average of segment OSRM distance
- population OSRM estimated distance is higher than the actual distance from source to destination warehouse.

## **EDA Results**

- we can observe that Mumbai Maharashtra ,Delhi , Gurgaon(Haryana),Bengaluru Karnataka ,Hyderabad Telangana, Chennai Tamil Nadu, Ahmedabad-Gujarat, Pune Maharashtra, Chandigarh Chandigarh and Kolkata West Bengal are some cities have highest amount of trips happening states with in the city.
- If we talk about , not having equal source and destination states , source and destination cities having highest number of trips in between are : Delhi TO Gurgaon , Gurgaon TO Bengaluru , Bhiwandi/Mumbai TO Pune Maharashtra , Sonipat TO Gurgaon,Haryana
- It is also been observed that lots of deliveries are happening to airports like : Chennai to MAA chennai international Airport , Pune to Pune Airport (PNQ),Kolkata to CCU West Bengal Kolkata International Airport , Bengaluru to BLR-Bengaluru International Airport etc.
- From Bar charts , and calculated tables in analysis , we can observe that highest trips are happening is with in the particular cities, in terms of average distance between destinations , we can observe Guwahati to Mumbai , Bengaluru to Chandigarh ,Bengaluru to Delhi , Bengaluru to Gurgaon are the longest routes.

## **Recommendations**

- As per analysis, It is recommended to use Carting (small vehicles) for delivery within the city in order to reduce the delivery time, and Heavy trucks for long distance trips or heavy load. based on this , we can optimize the delivery time as well as increase the revenue as per requirements.
- Increasing the connectivity in tier 2 and tier 3 cities along with professional tie-ups with several e-commerce giants can increase the revenue as well as the reputation on connectivity across borders.
- We can work on optimizing the scanning time on both ends which is start scanning time and end scanning time so that the delivery time can be equated to the OSRM estimated delivery time.
- Revisit information fed to routing engine for trip planning. Check for discrepancies with transporters, if the routing engine is configured for optimum results.

- North, South and West Zones comidors have significant traffic of orders. But, we have a smaller presence in Central, Eastern and North-Eastern zone. However it would be difficult to conclude this, by looking at just 2 months data. It is worth investigating and increasing our presence in these regions.
- From state point of view, we have heavy traffic in Mahrashta followed by Karnataka. This is a good indicator that we need to plan for resources on ground in these 2 states on priority. Especially, during festive seasons.