

## ▼ Business Problem Statement

The company wants to understand and process the data coming out of data engineering pipelines: Clean, sanitize and manipulate data to get useful features out of raw fields Make sense out of the raw data and help the data science team to build forecasting models on it.

The aim is to find if there are significant differences among variables like expected delivery time and actual delivery time. It is crucial for business because if the algorithm that predicts time is not working correctly, then customers will get wrong estimation and if the customers will not get delivery on time then, they will be unsatisfied and raise customer service complaints. That may lead to increase operational overload, decrease revenue, decrease company reputation.

## ▼ Exploratory Data Analysis

```
1 import numpy as np
2 import pandas as pd
3 from scipy import stats
4
5 import matplotlib.pyplot as plt
6 import seaborn as sns
```

```
1 df = pd.read_csv('/content/delhivery_data.txt')
```

```
1 pd.set_option('max_columns',100)
2 pd.set_option('min_rows', 8)
3 df
```

		data	trip_creation_time	route_schedule_uuid	route_type	trip_uuid	source_center	
0	training		2018-09-20 02:35:36.476840	thanos::sroute:eb7bfc78- b351-4c0e-a951- fa3d5c3...	Carting	153741093647649320	IND388121AAA	Anand
1	training		2018-09-20 02:35:36.476840	thanos::sroute:eb7bfc78- b351-4c0e-a951- fa3d5c3...	Carting	153741093647649320	IND388121AAA	Anand
2	training		2018-09-20 02:35:36.476840	thanos::sroute:eb7bfc78- b351-4c0e-a951- fa3d5c3...	Carting	153741093647649320	IND388121AAA	Anand
3	training		2018-09-20 02:35:36.476840	thanos::sroute:eb7bfc78- b351-4c0e-a951- fa3d5c3...	Carting	153741093647649320	IND388121AAA	Anand
...	...	...	...	...	...	...	...	...
144863	training		2018-09-20 16:24:28.436231	thanos::sroute:f0569d2f- 4e20-4c31-8542- 67b86d5...	Carting	153746066843555182	IND131028AAB	Sor
144864	training		2018-09-20 16:24:28.436231	thanos::sroute:f0569d2f- 4e20-4c31-8542- 67b86d5...	Carting	153746066843555182	IND131028AAB	Sor
144865	training		2018-09-20 16:24:28.436231	thanos::sroute:f0569d2f- 4e20-4c31-8542- 67b86d5...	Carting	153746066843555182	IND131028AAB	Sor
144866	training		2018-09-20 16:24:28.436231	thanos::sroute:f0569d2f- 4e20-4c31-8542- 67b86d5...	Carting	153746066843555182	IND131028AAB	Sor

144867 rows × 24 columns

```
1 df.shape
2 # 144867 rows, 24 columns
```

(144867, 24)

```
1 # Drop unnecessary/Unknown field columns
2 unknown = df.iloc[:,df.columns.str.contains('factor|cutoff')].columns
3 for i in unknown:
4     df.drop(i,axis=1,inplace=True)
```

```
1 df_na = pd.DataFrame(df.isna().sum())
2 df_na['percent'] = df.isna().sum() *100/len(df)
```

```

3 df_na['percent'] = df_na['percent'].round(3)
4 df_na
5 # Minor rows with missing value, we can drop those rows.

```

	0	percent
<b>data</b>	0	0.000
<b>trip_creation_time</b>	0	0.000
<b>route_schedule_uuid</b>	0	0.000
<b>route_type</b>	0	0.000
<b>trip_uuid</b>	0	0.000
<b>source_center</b>	0	0.000
<b>source_name</b>	293	0.202
<b>destination_center</b>	0	0.000
<b>destination_name</b>	261	0.180
<b>od_start_time</b>	0	0.000
<b>od_end_time</b>	0	0.000
<b>start_scan_to_end_scan</b>	0	0.000
<b>actual_distance_to_destination</b>	0	0.000
<b>actual_time</b>	0	0.000
<b>osrm_time</b>	0	0.000
<b>osrm_distance</b>	0	0.000
<b>segment_actual_time</b>	0	0.000
<b>segment_osrm_time</b>	0	0.000
<b>segment_osrm_distance</b>	0	0.000

```

1 df.dropna(how='any', inplace=True)
2 # removed rows with null values

```

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

```

<class 'pandas.core.frame.DataFrame'>
Int64Index: 144316 entries, 0 to 144866
Data columns (total 19 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   data             144316 non-null   object  
 1   trip_creation_time 144316 non-null   object  
 2   route_schedule_uuid 144316 non-null   object  
 3   route_type        144316 non-null   object  
 4   trip_uuid          144316 non-null   object  
 5   source_center      144316 non-null   object  
 6   source_name        144316 non-null   object  
 7   destination_center 144316 non-null   object  
 8   destination_name   144316 non-null   object  
 9   od_start_time      144316 non-null   object  
 10  od_end_time        144316 non-null   object  
 11  start_scan_to_end_scan 144316 non-null   float64
 12  actual_distance_to_destination 144316 non-null   float64
 13  actual_time        144316 non-null   float64
 14  osrm_time          144316 non-null   float64
 15  osrm_distance      144316 non-null   float64
 16  segment_actual_time 144316 non-null   float64
 17  segment_osrm_time  144316 non-null   float64
 18  segment_osrm_distance 144316 non-null   float64
dtypes: float64(8), object(11)
memory usage: 22.0+ MB

```

```

1 df['od_start_time'] = pd.to_datetime(df['od_start_time'])
2 df['od_end_time'] = pd.to_datetime(df['od_end_time'])
3 # Converting the data type of columns to appropriate data type

```

```

1 df.nunique()
2 # Only few columns are categorical, mostly data is numeric

```

<b>data</b>	2
<b>trip_creation_time</b>	14787
<b>route_schedule_uuid</b>	1497
<b>route_type</b>	2
<b>trip_uuid</b>	14787
<b>source_center</b>	1496

```

source_name           1496
destination_center    1466
destination_name       1466
od_start_time        26223
od_end_time          26223
start_scan_to_end_scan   1914
actual_distance_to_destination 143965
actual_time            3182
osrm_time              1531
osrm_distance         137544
segment_actual_time     746
segment_osrm_time      214
segment_osrm_distance  113497
dtype: int64

```

```

1 #Checking how the data is spread in categoric columns
2 df2=df.copy()
3 df2
4 categ_cols = ['route_type']
5 cat_count = df2[categ_cols].melt().groupby(['variable', 'value'])[['value']].size().reset_index(name='counts')
6 s = df2[categ_cols].melt().variable.value_counts()
7 cat_count['Percent'] = cat_count['counts'].div(cat_count['variable'].map(s)).mul(100).round().astype('int')
8 cat_count.groupby(['variable', 'value']).first()
9 # FTL Transportation used 69%
10 # Carting Transportation 31%

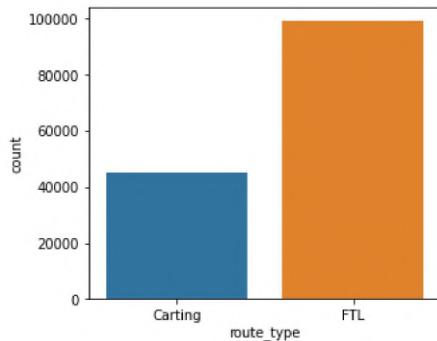
```

		counts	Percent
variable	value		
route_type	Carting	45184	31
	FTL	99132	69

```

1 plt.figure(figsize = [16,4])
2 cat_cols = ['route_type']
3 for i in range (len(cat_cols)):
4     plt.subplot(1, 3, i+1)
5     sns.countplot(data=df, x=cat_cols[i])

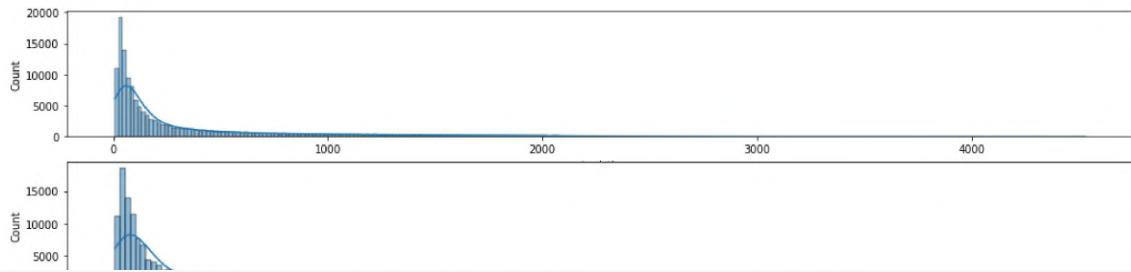
```



```

1 plt.figure(figsize = [18,10])
2 num_cols = ['actual_time','osrm_time','actual_distance_to_destination','osrm_distance']
3 for i in range (len(num_cols)):
4     plt.subplot(len(num_cols),1, i+1)
5     sns.histplot(data=df, x=num_cols[i], kde=True)

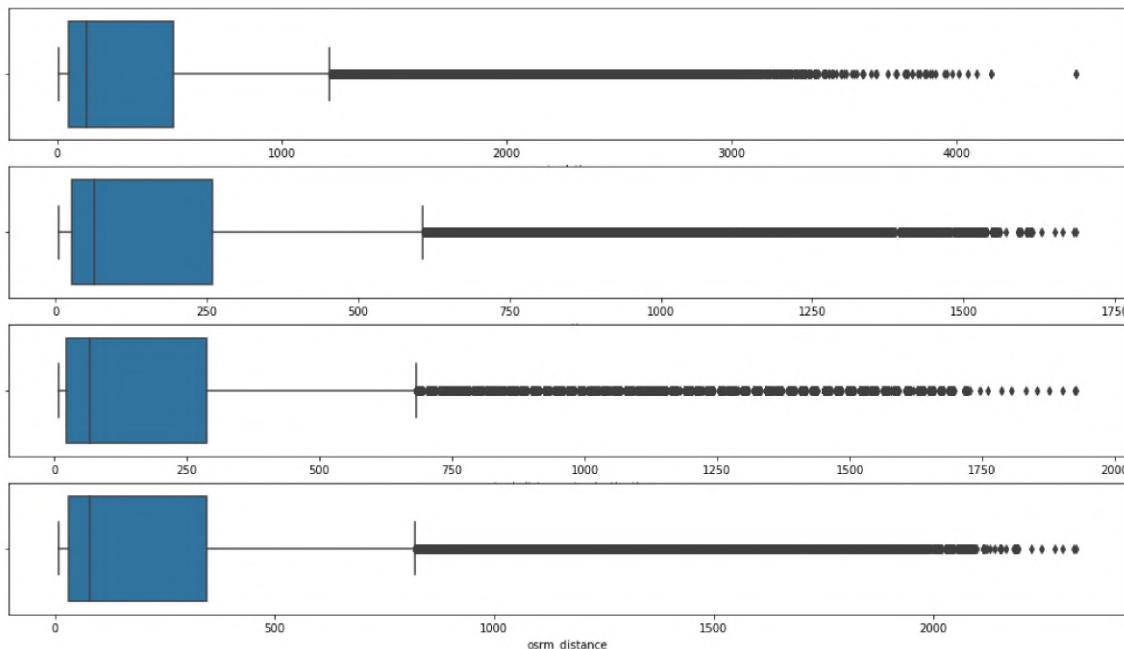
```



```

1 plt.figure(figsize = [18,10])
2 num_cols = ['actual_time','osrm_time','actual_distance_to_destination','osrm_distance']
3 for i in range (len(num_cols)):
4     plt.subplot(len(num_cols),1, i+1)
5     sns.boxplot(data=df, x=num_cols[i])

```



```

1 df.describe()
2
3 # Data contains extreame outliers
4 # Variables showing are highly right skew distribution
5 # segment_actual_time showing negetive value that not making sense
6 # Data cleaning is needed
7 # Data contains similer segment vaulues, need to group them

```

	start_scan_to_end_scan	actual_distance_to_destination	actual_time	osrm_time	osrm_distance	segm
count	130868.000000	130868.000000	130868.000000	130868.000000	130868.000000	
mean	941.958057	231.120919	404.065012	209.825603	279.976641	
std	1031.953320	342.810982	591.574488	305.929392	418.394128	
min	20.000000	9.000045	9.000000	6.000000	9.008200	
25%	151.000000	23.124073	48.000000	26.000000	28.792675	
50%	426.000000	63.811633	119.000000	61.000000	75.009050	
75%	1597.000000	286.273519	488.000000	251.000000	335.047625	
max	4535.000000	1924.221719	4532.000000	1682.000000	2319.934400	

## ▼ groupby and aggregations

```

1 # Grouping the data at journey level, Creating data with unique index value as segment_key.
2 df['segment_key'] = df['trip_uuid'] + df['source_center'] + df['destination_center']
3 segment_cols = ['segment_actual_time','segment_osrm_distance','segment_osrm_time']

```

```

4 for col in segment_cols:
5     df[col+'_sum'] = df.groupby('segment_key')[col].cumsum()
6 df[[col + '_sum' for col in segment_cols]].head()

```

	segment_actual_time_sum	segment_osrm_distance_sum	segment_osrm_time_sum
0	14.0	11.9653	11.0
1	24.0	21.7243	20.0
2	40.0	32.5395	27.0
3	61.0	45.5619	39.0
4	67.0	49.4772	44.0

```

1 # Aggregating at sub-journey level
2 create_segment_dict = {
3     'data' : 'first',
4     'trip_creation_time' : 'first',
5     'route_schedule_uuid' : 'first',
6     'route_type' : 'first',
7     'trip_uuid' : 'first',
8
9     'source_center' : 'first',
10    'source_name' : 'first',
11
12    'destination_center' : 'last',
13    'destination_name' : 'last',
14
15    'od_start_time' : 'first',
16    'od_end_time' : 'first',
17    'start_scan_to_end_scan' : 'first',
18
19    'actual_distance_to_destination' : 'last',
20    'actual_time' : 'last',
21
22    'osrm_time' : 'sum',
23    'osrm_distance' : 'sum',
24
25    'segment_actual_time_sum' : 'sum',
26    'segment_osrm_distance_sum' : 'sum',
27    'segment_osrm_time_sum' : 'sum'
28 }

```

```

1 segment = df.groupby('segment_key').agg(create_segment_dict).reset_index()
2 segment = segment.sort_values(by=['segment_key','od_end_time'],ascending=True).reset_index()

```

## ▼ feature creation and merging of rows

```

1 # Feature creation: Calculating time taken between od_start_time and od_end_time, keep it as feature.
2 segment['od_time_diff_hour'] = (segment['od_end_time'] - segment['od_start_time']).dt.total_seconds()/60

```

```

1 segment
2 # od_time_diff_hour is matching with start_scan_to_end_scan

```

index		segment_key	data	trip_creation_time	route_schedule_uuid
0	0	153671041653548748IND209304AAIND00000ACB	trip-training	2018-09-12 00:00:16.535741	thanos::sroute:d7c989ba-a29b-4a0b-b2f4-288cdc6...
1	1	153671041653548748IND462022AAIND209304AAA	trip-training	2018-09-12 00:00:16.535741	thanos::sroute:d7c989ba-a29b-4a0b-b2f4-288cdc6...
2	2	153671041653548748IND5612022AAIND662101AAA	trip-training	2018-09-12 00:00:22.896130	thanos::sroute:3a1b0ab2-bb0b-4c53-8c59-

```

1 segment.nunique()
2 #now we have the data without duplicates, means we have unique value for ever row as segment_key.

```

```

index          26222
segment_key    26222
data            2
trip_creation_time 14787
route_schedule_uuid 1497
route_type      2
trip_uuid        14787
source_center    1496
source_name       1496
destination_center 1466
destination_name 1466
od_start_time   26222
od_end_time     26222
start_scan_to_end_scan 1914
actual_distance_to_destination 26193
actual_time      1657
osrm_time        2723
osrm_distance    26089
segment_actual_time_sum 3839
segment_osrm_distance_sum 26102
segment_osrm_time_sum 2964
od_time_diff_hour 26222
dtype: int64

```

```

1 create_trip_dict = {
2     'data' : 'first',
3     'trip_creation_time' : 'first',
4     'route_schedule_uuid' : 'first',
5     'route_type' : 'first',
6     'trip_uuid' : 'first',
7
8     'source_center' : 'first',
9     'source_name' : 'first',
10
11    'destination_center' : 'last',
12    'destination_name' : 'last',
13
14    'start_scan_to_end_scan' : 'sum',
15    'actual_distance_to_destination' : 'sum',
16    'actual_time' : 'sum',
17    'osrm_time' : 'sum',
18    'osrm_distance' : 'sum',
19
20    'segment_actual_time_sum' : 'sum',
21    'segment_osrm_distance_sum' : 'sum',
22    'segment_osrm_time_sum' : 'sum'
23 }

```

```

1 trip = segment.groupby('trip_uuid').agg(create_trip_dict).reset_index(drop=True)
2 trip

```

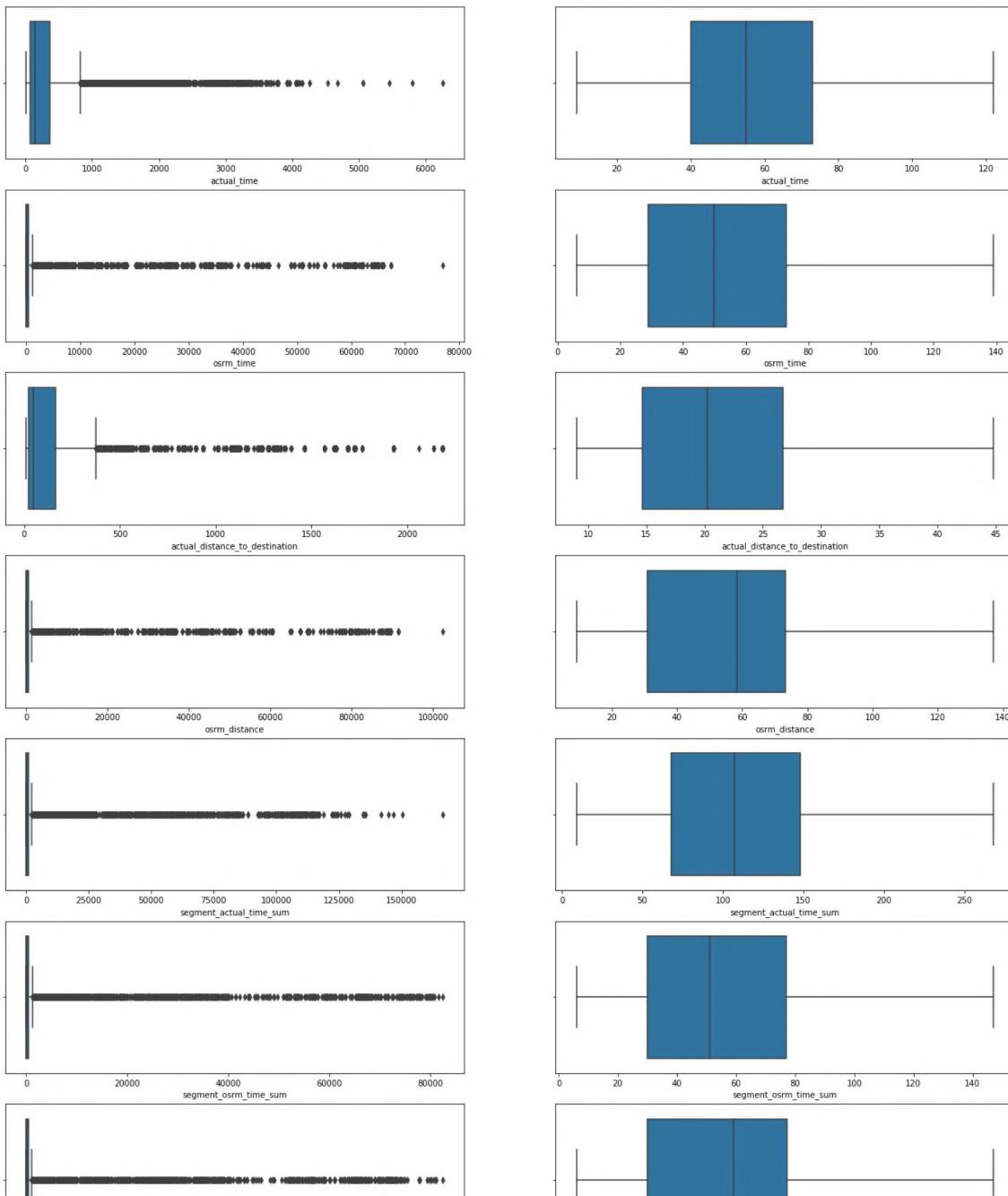
	data	trip_creation_time	route_schedule_uuid	route_type		trip_uuid	source_center
0	training	2018-09-12 00:00:16.535741	thanos::sroute:d7c989ba-a29b-4a0b-b2f4-288cdc6...	FTL	trip-	153671041653548748	IND209304AAA Ka
1	training	2018-09-12 00:00:22.886430	thanos::sroute:3a1b0ab2-bb0b-4c53-8c59-eb2a2c0...	Carting	trip-	153671042288605164	IND561203AAB Doddak
2	training	2018-09-12 00:00:33.691250	thanos::sroute:de5e208e-7641-45e6-8100-4d9fb1e...	FTL	trip-	153671043369099517	IND000000ACB Gur
3	training	2018-09-12 00:01:00.113710	thanos::sroute:f0176492-a679-4597-8332-bbd1c7f...	Carting	trip-	153671046011330457	IND400072AAB
...	...	...	...	...	...	...	...

## ▼ Outlier treatment

```

1 trip_o = trip.copy()
2
3 def find_outliers_IQR(col):
4     q1=col.quantile(0.25)
5     q3=col.quantile(0.75)
6     IQR=q3-q1
7     outliers = trip[((col<(q1-1.5*IQR)) | (col>(q3+1.5*IQR)))]
8     return outliers
9
10 cols = ['actual_time', 'osrm_time','actual_distance_to_destination','osrm_distance','segment_actual_time_sum','segment_osrm_time_sum',
11
12 n=1
13 while n!=0:
14     n=0
15     for x in cols:
16         outliers = find_outliers_IQR(trip[x]).index
17         trip.drop(outliers,inplace=True)
18         n+=len(outliers)
19
20
21 fig, axis = plt.subplots(nrows=len(cols), ncols=2, figsize=(22, len(cols)*4))
22 for i in range (len(cols)):
23     for j in ([0,1]):
24         if j==0:
25             sns.boxplot(data=trip_o, x=cols[i], ax=axis[i, j])
26         else:
27             sns.boxplot(data=trip, x=cols[i], ax=axis[i, j])
28 # Left side plots: boxplot distribution before removing outliers
29 # Right side plots: boxplot distribution after removing outliers

```



```

1 trip_o.describe()
2 # before outlier treatment there was high difference between median and mean values.

```

	start_scan_to_end_scan	actual_distance_to_destination	actual_time	osrm_time	osrm_distance	segmer
count	14787.000000	14787.000000	14787.000000	14787.000000	14787.000000	14787.000000
mean	529.429025	164.090196	356.306012	2092.831406	2786.867034	
std	658.254936	305.502982	561.517936	7964.769556	10769.741065	
min	23.000000	9.002461	9.000000	6.000000	9.072900	
25%	149.000000	22.777099	67.000000	61.000000	65.635550	
50%	279.000000	48.287894	148.000000	167.000000	172.886300	
75%	632.000000	163.591258	367.000000	510.500000	603.955400	
max	7898.000000	2186.531787	6265.000000	76953.000000	102415.868000	

```

1 trip.describe()
2 # after outlier treatment there the mean values is closer to median.

```

	start_scan_to_end_scan	actual_distance_to_destination	actual_time	osrm_time	osrm_distance	segment_
<b>count</b>	5160.000000	5160.000000	5160.000000	5160.000000	5160.000000	
<b>mean</b>	155.022481	20.817280	57.434690	54.112209	57.274155	
<b>std</b>	121.042528	8.006889	23.920582	29.317440	29.436978	
<b>min</b>	23.000000	9.002461	9.000000	6.000000	9.072900	
<b>25%</b>	89.000000	14.641301	40.000000	29.000000	30.853800	
<b>50%</b>	127.000000	20.266599	55.000000	50.000000	58.327150	
<b>75%</b>	182.000000	26.713161	73.000000	73.000000	73.351800	

```
1 len(trip_o),len(trip)

(14787, 5160)
```

## ▼ Hypothesis testing Comparisons

```
1 # hypothesis testing/ visual analysis between actual_time aggregated value and OSRM time aggregated value
2
3 # Checking: Does segment_actual_time is similer as segment_osrm_time?
4
5 from scipy.stats import ttest_ind
6 null_hypothesis = 'mean of actual_time is not higher than mean of osrm_time'
7 alternative_hypothesis = 'mean of actual_time is higher than mean of osrm_time'
8
9 sample1 = trip['actual_time']
10 sample2 = trip['osrm_time']
11 t_stat, p_value = ttest_ind(sample1, sample2, equal_var=False, alternative='greater')
12 print(t_stat, p_value)
13
14 if(p_value < 0.05):
15     print('Since, p-value < 0.05, the null hypothesis is rejected')
16     print(alternative_hypothesis)
17 else:
18     print('Since p-value > 0.05, we fail to reject null hypothesis')
19     print(null_hypothesis)
20
21 # conclusion: mean of actual_time is higher than mean of osrm_time
```

6.307544212335001 1.4779513885349015e-10  
 Since, p-value < 0.05, the null hypothesis is rejected  
 mean of actual\_time is higher than mean of osrm\_time

```
1 # hypothesis testing/ visual analysis between actual_time aggregated value and segment actual time aggregated value
2
3 # Checking: Does actual_time is similer as segment_osrm_time?
4
5 from scipy.stats import ttest_ind
6 null_hypothesis = 'mean of actual_time is similer as segment_actual_time'
7 alternative_hypothesis = 'mean of actual_time is different than mean of segment_osrm_time'
8
9 sample1 = trip['actual_time']
10 sample2 = trip['segment_actual_time_sum']
11 t_stat, p_value = ttest_ind(sample1, sample2)
12 print(t_stat, p_value)
13
14 if(p_value < 0.05):
15     print('Since, p-value < 0.05, the null hypothesis is rejected')
16     print(alternative_hypothesis)
17 else:
18     print('Since p-value > 0.05, we fail to reject null hypothesis')
19     print(null_hypothesis)
20
21 # conclusion: mean of actual_time is different than mean of segment_osrm_time
```

-65.79091499324902 0.0  
 Since, p-value < 0.05, the null hypothesis is rejected  
 mean of actual\_time is different than mean of segment\_osrm\_time

```
1 # Checking: Does osrm_distance is similer as segment_osrm_distance_sum?
2
3 from scipy.stats import ttest_ind
4 null_hypothesis = 'mean of osrm_distance is similer as mean of segment_osrm_distance_sum'
5 alternative_hypothesis = 'mean of osrm_distance is higher than mean of segment_osrm_distance_sum'
6
```

```

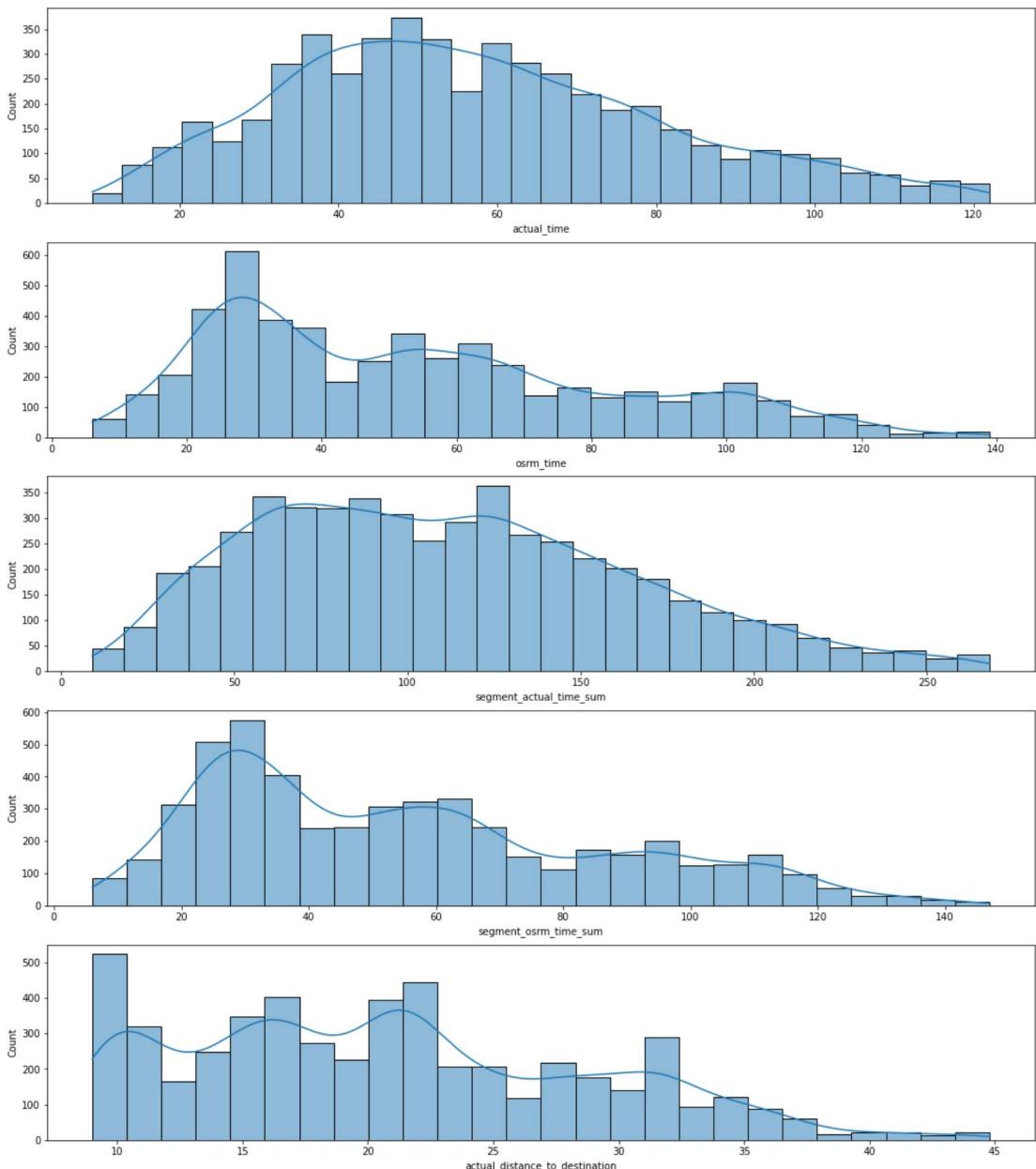
7 sample1 = trip['osrm_distance']
8 sample2 = trip['segment_osrm_distance_sum']
9 t_stat, p_value = ttest_ind(sample1, sample2, equal_var=False, alternative='greater')
10 print(t_stat, p_value)
11
12 if(p_value < 0.05):
13     print('Since, p-value < 0.05, the null hypothesis is rejected')
14     print(alternative_hypothesis)
15 else:
16     print('Since p-value > 0.05, we fail to reject null hypothesis')
17     print(null_hypothesis)
18
19 # conclusion: mean of osrm_distance is similer as mean of segment_osrm_distance_sum
-3.4156110488999936 0.9996805751383507
Since p-value > 0.05, we fail to reject null hypothesis
mean of osrm_distance is similer as mean of segment_osrm_distance_sum

1 num_cols = ['actual_time','osrm_time','segment_actual_time_sum','segment_osrm_time_sum','actual_distance_to_destination','osrm_distance']
2 for i in (num_cols):
3     stat, p_value = stats.shapiro(sample1)
4     if(p_value < 0.05):
5         print(i, ": sample is not normally distributed, do non parametric test")
6     else:
7         print(i, ": sample is normally distributed, can do parametric test")

actual_time : sample is not normally distributed, do non parametric test
osrm_time : sample is not normally distributed, do non parametric test
segment_actual_time_sum : sample is not normally distributed, do non parametric test
segment_osrm_time_sum : sample is not normally distributed, do non parametric test
actual_distance_to_destination : sample is not normally distributed, do non parametric test
osrm_distance : sample is not normally distributed, do non parametric test
segment_osrm_distance_sum : sample is not normally distributed, do non parametric test
/usr/local/lib/python3.8/dist-packages/scipy/stats/morestats.py:1760: UserWarning: p-value may not be accurate for N > 5000.
warnings.warn("p-value may not be accurate for N > 5000.")

1 plt.figure(figsize = [18,30])
2 num_cols = ['actual_time','osrm_time','segment_actual_time_sum','segment_osrm_time_sum','actual_distance_to_destination','osrm_distance']
3 for i in range (len(num_cols)):
4     plt.subplot(len(num_cols),1, i+1)
5     sns.histplot(data=trip, x=num_cols[i], kde=True)
6 # Distribution are not normally distributed, we will do non parametric tests

```



```

1 # H0: mean of both samples are similer
2 # Ha: means of both samples are different
3 from scipy.stats import mannwhitneyu
4
5 sample1 = trip['actual_time']
6 sample2 = trip['osrm_time']
7 # perform mann whitney test
8 stat, p_value = mannwhitneyu(sample1, sample2)
9 print('Statistics=%.2f, p=%.2f' % (stat, p_value))
10 # Level of significance
11 alpha = 0.05
12 # conclusion
13 if p_value < alpha:
14     print('Reject Null Hypothesis (Significant difference between two samples)')
15 else:
16     print('Do not Reject Null Hypothesis (No significant difference between two samples)')

Statistics=14829121.00, p=0.00
Reject Null Hypothesis (Significant difference between two samples)

```

```

1 # H0: mean of both samples are similer
2 # Ha: means of both samples are different
3 from scipy.stats import mannwhitneyu
4
5 sample1 = trip['osrm_distance']
6 sample2 = trip['segment_osrm_distance_sum']
7 # perform mann whitney test
8 stat, p_value = mannwhitneyu(sample1, sample2)
9 print('Statistics=%.2f, p=%.2f' % (stat, p_value))
10 # Level of significance
11 alpha = 0.05
12 # conclusion

```

```

13 if p_value < alpha:
14     print('Reject Null Hypothesis (Significant difference between two samples)')
15 else:
16     print('Do not Reject Null Hypothesis (No significant difference between two samples)')

Statistics=12813705.50, p=0.00
Reject Null Hypothesis (Significant difference between two samples)

```

## ▼ Standardization, Normalization

```

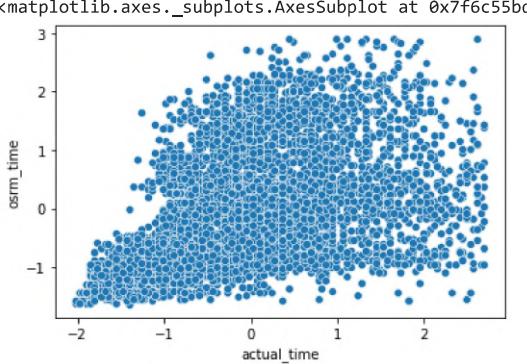
1 df_ao = trip[["actual_time", "osrm_time"]]

1 from sklearn.preprocessing import StandardScaler, MinMaxScaler
2 df_ao_ss = StandardScaler().fit_transform(df_ao) # ss--> standard scaler z-score

1 df_ao_ss = pd.DataFrame(df_ao_ss, columns=["actual_time", "osrm_time"])

1 sns.scatterplot(x=df_ao_ss["actual_time"], y=df_ao_ss["osrm_time"])

```

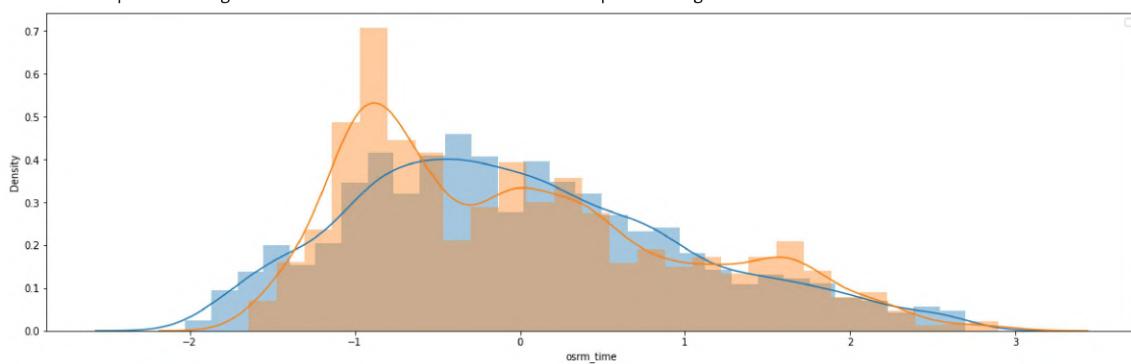


```

1 plt.figure(figsize = (20,6))
2
3 sns.distplot(df_ao_ss['actual_time'])
4 sns.distplot(df_ao_ss['osrm_time'])
5
6 plt.legend()
7 plt.show()

/usr/local/lib/python3.8/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is a deprecated
    warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.8/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is a deprecated
    warnings.warn(msg, FutureWarning)
WARNING:matplotlib.legend:No handles with labels found to put in legend.

```



```

1 # hypothesis testing/ visual analysis between actual_time aggregated value and OSRM time aggregated value
2
3 # Checking: Does segment_actual_time is similer as segment_osrm_time?
4
5 from scipy.stats import ttest_ind
6 null_hypothesis = 'mean of actual_time is similer to osrm_time'
7 alternative_hypothesis = 'mean of actual_time is different than osrm_time'

```

```

8
9 sample1 = df_ao_ss['actual_time']
10 sample2 = df_ao_ss['osrm_time']
11 t_stat, p_value = ttest_ind(sample1, sample2)
12 print(t_stat, p_value)
13
14 if(p_value < 0.05):
15     print('Since, p-value < 0.05, the null hypothesis is rejected')
16     print(alternative_hypothesis)
17 else:
18     print('Since p-value > 0.05, we fail to reject null hypothesis')
19     print(null_hypothesis)
20
21 # conclusion: mean of actual_time is similer to osrm_time (with following the standardization approach)
-4.301138070642346e-15 0.9999999999999966
Since p-value > 0.05, we fail to reject null hypothesis
mean of actual_time is similer to osrm_time

```

```

1 df_ao_ss.mean()

actual_time    -8.950635e-17
osrm_time      -4.819573e-18
dtype: float64

```

```

1 df_ao_mm.mean()

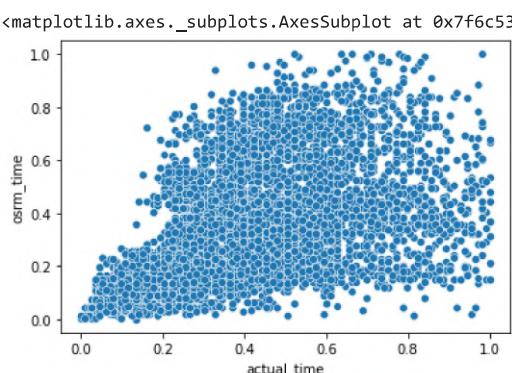
actual_time    0.428626
osrm_time      0.361746
dtype: float64

```

```

1 df_ao_mm = MinMaxScaler().fit_transform(df_ao)
2 df_ao_mm = pd.DataFrame(df_ao_mm, columns=["actual_time", "osrm_time"])
3 sns.scatterplot(x=df_ao_mm["actual_time"], y=df_ao_mm["osrm_time"])

```



```

1 plt.figure(figsize = (20,6))
2
3 sns.distplot(df_ao_mm['actual_time'])
4 sns.distplot(df_ao_mm['osrm_time'])
5
6 plt.legend()
7 plt.show()

```

```
/usr/local/lib/python3.8/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is a deprecated
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.8/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is a deprecated
warnings.warn(msg, FutureWarning)
```

## Handling categorical values

```
1 trip.nunique()

data           2
trip_creation_time    5160
route_schedule_uuid   746
route_type          2
trip_uuid          5160
source_center       396
source_name         396
destination_center  415
destination_name    415
start_scan_to_end_scan  500
actual_distance_to_destination  5151
actual_time        114
osrm_time          134
osrm_distance      5140
segment_actual_time_sum  260
segment_osrm_distance_sum  5142
segment_osrm_time_sum   141
dtype: int64
```

```
1 from sklearn.preprocessing import LabelEncoder
2 label_encoder = LabelEncoder()
3 trip[col] = label_encoder.fit_transform(trip['route_type'])
4 trip[col].value_counts()

0    4907
1     253
Name: route_type, dtype: int64
```

```
1 trip['data'].value_counts()

training    3599
test        1561
Name: data, dtype: int64
```

```
1 from sklearn.preprocessing import LabelEncoder
2 label_encoder = LabelEncoder()
3 trip[col] = label_encoder.fit_transform(trip['data'])
4 trip[col].value_counts()

1    3599
0    1561
Name: route_type, dtype: int64
```

## Handling missing values

```
1 from sklearn.impute import SimpleImputer
2 dm = pd.read_csv('/content/delhivery_data.txt')
3 dm.head(3)
4
5 # Because data contains very few missing values, so we removed them.
6 # If we would get data with many missing value, then we would imputation approach:
```

	data	trip_creation_time	route_schedule_uuid	route_type	trip_uuid	source_center	source_name
0	training	2018-09-20 02:35:36.476840	thanos::sroute:eb7bfc78- b351-4c0e-a951- fa3d5c3...	Carting	trip- 153741093647649320	IND388121AAA	Anand_VUN@ (
1	training	2018-09-20 02:35:36.476840	thanos::sroute:eb7bfc78- b351-4c0e-a951- fa3d5c3...	Carting	trip- 153741093647649320	IND388121AAA	Anand_VUN@ (
2	training	2018-09-20 02:35:36.476840	thanos::sroute:eb7bfc78- b351-4c0e-a951- fa3d5c3...	Carting	trip- 153741093647649320	IND388121AAA	Anand_VUN@ (



```
1 dm.isna().sum()

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

```
1 dm['source_name'] = SimpleImputer(strategy="most_frequent").fit_transform(dm[['source_name']])

1 dm['destination_name'] = SimpleImputer(strategy="most_frequent").fit_transform(dm[['destination_name']])
```

```
1 dm.isna().sum()

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

## ▼ Columns split

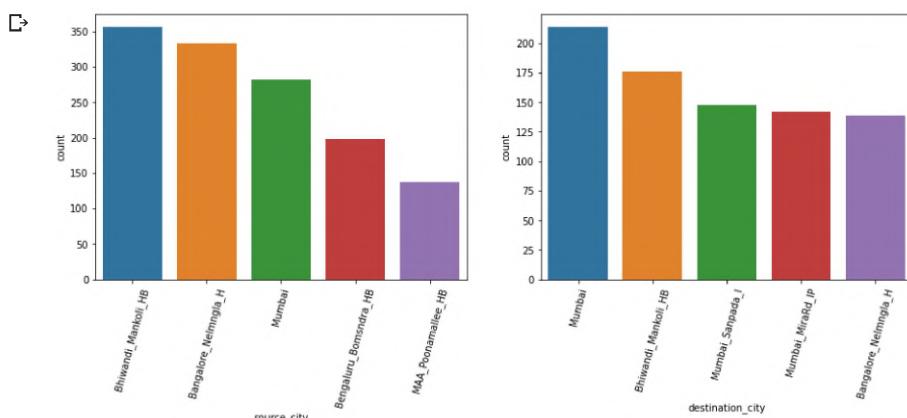
```
1 ds = trip[['destination_name']].copy()
2
3 new = trip['source_name'].str.split(" ", n = 1, expand = True)
4 ds['source_city']= new[0]
5 ds['source_state']= new[1].str[1:-1]
6
7 new = trip['destination_name'].str.split(" ", n = 1, expand = True)
8 ds['destination_city']= new[0]
9 ds['destination_state']= new[1].str[1:-1]
10
11 ds['Corridor'] = ds['source_city']+ To "+ds['destination_city']"
12
13 ds
```

	destination_name	source_city	source_state	destination_city	destination_state
3	Mumbai_MiraRd_IP (Maharashtra)	Mumbai	ub (Maharashtra)	Mumbai_MiraRd_IP	Maharashtra
5	Chennai_Poonamallee (Tamil Nadu)	Chennai_Poonamallee	Tamil Nadu	Chennai_Poonamallee	Tamil Nadu
6	Chennai_Vandalur_Dc (Tamil Nadu)	Chennai_Chrompet_DPC	Tamil Nadu	Chennai_Vandalur_Dc	Tamil Nadu
7	HBR Layout PC (Karnataka)	HBR	ayout PC (Karnataka)	HBR	ayout PC (Karnataka)
...	...	...	...	...	...
14779	Chennai_Thiruvl_DC (Tamil Nadu)	Chennai_Porur_DPC	Tamil Nadu	Chennai_Thiruvl_DC	Tamil Nadu

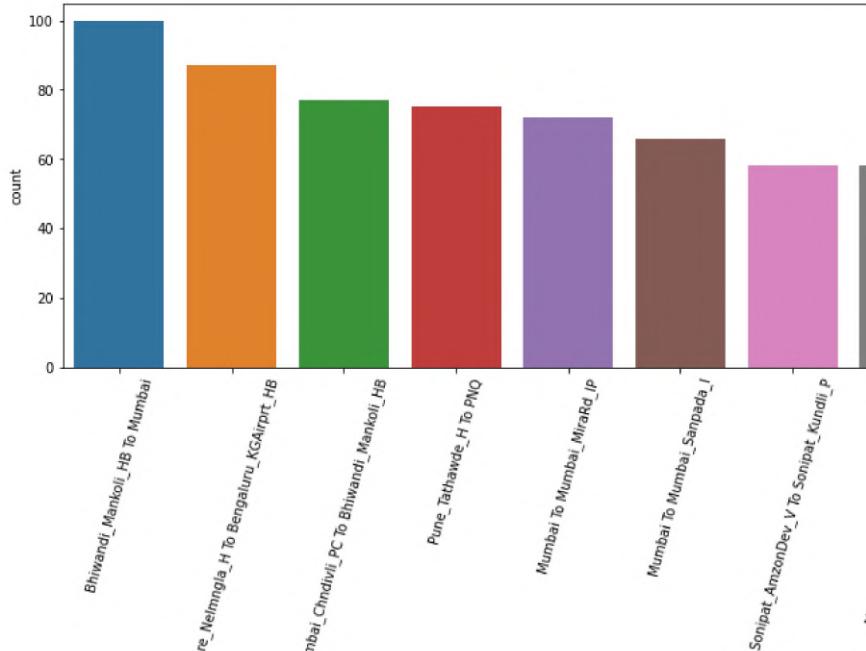
```
1 ds['Corridor'].value_counts()
2 # There are 658 routes
```

```
Bhiwandi_Mankoli_HB To Mumbai 100
Bangalore_Nelmgla_H To Bengaluru_KGAirprt_HB 87
Mumbai_Chndivli_PC To Bhiwandi_Mankoli_HB 77
Pune_Tathawde_H To PNQ 75
...
Delhi_Rohini_DPC To Delhi_Barwala 1
Wardha_RamaNgr_D To Deoli_Central_DPP_2 1
Nadiad_DC To Nadiad_DC 1
Janakpuri To Delhi_Nangli_IP 1
Name: Corridor, Length: 658, dtype: int64
```

```
1 plt.figure(figsize=(15,5))
2 plt.subplot(1,2,1)
3 sns.countplot(data= ds, x='source_city', order=ds['source_city'].value_counts().nlargest(5).index)
4 plt.xticks(rotation = 75)
5
6 plt.subplot(1,2,2)
7 sns.countplot(data= ds, x='destination_city', order=ds['destination_city'].value_counts().nlargest(5).index)
8 plt.xticks(rotation = 75)
9 plt.show()
10
11 # most orders are coming from Bhindwandi_Mankoli_HB
12 # most orders are going to Mumbai
13 # Left plot: top 5 cities acting as source point
14 # right plot: top 5 cities acting as destination point
```



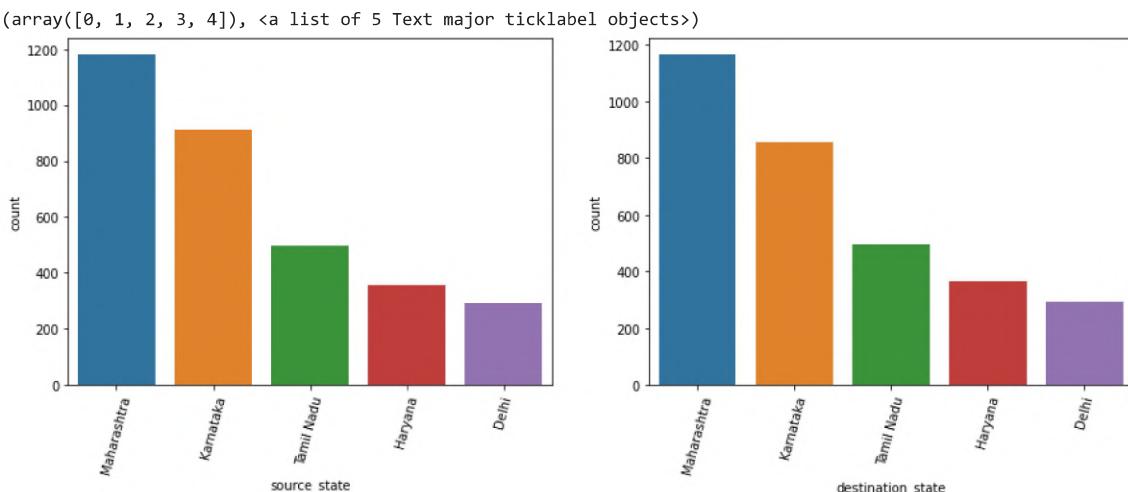
```
1 plt.figure(figsize=(15,5))
2 sns.countplot(data= ds, x='Corridor', order=ds['Corridor'].value_counts().nlargest(10).index)
3 plt.xticks(rotation = 75)
4 plt.show()
5
6 # the busiest route is Bhindwandi_Mankoli_HB to Mumbai
7 # Top 10 busiest routes
```



```

1 plt.figure(figsize=(15,5))
2
3 plt.subplot(1,2,1)
4 sns.countplot(data= ds, x='source_state', order=ds['source_state'].value_counts().nlargest(5).index)
5 plt.xticks(rotation = 75)
6
7 plt.subplot(1,2,2)
8 sns.countplot(data= ds, x='destination_state', order=ds['destination_state'].value_counts().nlargest(5).index)
9 plt.xticks(rotation = 75)
10
11 # most orders are coming from Maharashtra state
12 # most orders are going to Maharashtra state
13 # Left plot: top 5 cities acting as source point
14 # right plot: top 5 cities acting as destination point

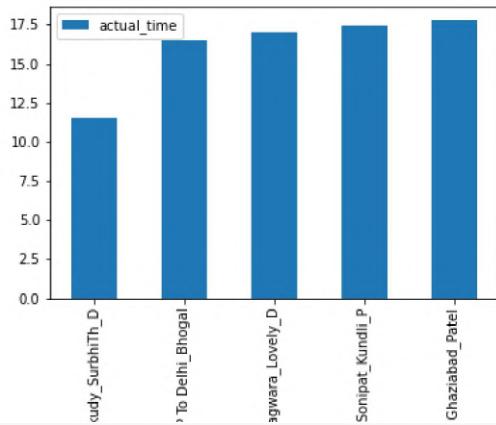
```



```

1 dn=pd.concat([trip,ds],axis=1)
2
3
4 dn.groupby('Corridor').agg({'actual_time':'mean'}).nsmallest(5,columns='actual_time').plot(kind='bar')
5 plt.show()
6
7 # trip between cities Angamaly to Chalakudy saw the least avg time for completion

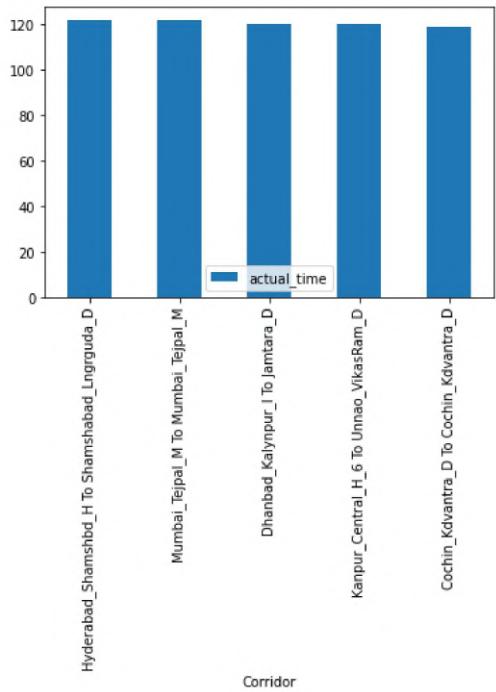
```



```

1 dn.groupby('Corridor').agg({'actual_time':'mean'}).nlargest(5,columns='actual_time').plot(kind='bar')
2 plt.show()
3
4 # trip between cities Hyderabad to Shamshabad saw the highest avg time for completion

```



```

1 ds.describe()
2 # 1179 orders are coming from state - Maharashtra

```

	destination_name	source_city	source_state	destination_city	destination_state	Corr
count	5160	5160	5160	5160	5160	5160
unique	415	391	42	411	44	
top	Mumbai Hub (Maharashtra)	Bhiwandi_Mankoli_HB	Maharashtra	Mumbai	Maharashtra	Bhiwandi_Manko To Mu
freq	212	357	1179	214	1164	

```

1 dn.describe()
2 # Average actual time for a trip is 57 Hour

```

```

route_type start_scan_to_end_scan actual_distance_to_destination actual_time osrm_time osrm_dist
count 5160.000000 5160.000000 5160.000000 5160.000000 5160.000000 5160.000000
mean 0.607191 155.022191 20.817280 57.421600 51.112200 57.271
1 trip_cr = df[['trip_creation_time']].copy()
2 trip_cr['trip_creation_time'] = pd.to_datetime(trip_cr['trip_creation_time'])
3 trip_cr['year'] = trip_cr['trip_creation_time'].dt.year
4 trip_cr['month'] = trip_cr['trip_creation_time'].dt.month
5 trip_cr['day'] = trip_cr['trip_creation_time'].dt.day
6 trip_cr
7 # Trip creation time variable is not making much sense

```

	trip_creation_time	year	month	day	edit
0	2018-09-20 02:35:36.476840	2018	9	20	
1	2018-09-20 02:35:36.476840	2018	9	20	
2	2018-09-20 02:35:36.476840	2018	9	20	
3	2018-09-20 02:35:36.476840	2018	9	20	
...	...	...	...	...	
144862	2018-09-20 16:24:28.436231	2018	9	20	
144863	2018-09-20 16:24:28.436231	2018	9	20	
144864	2018-09-20 16:24:28.436231	2018	9	20	
144865	2018-09-20 16:24:28.436231	2018	9	20	

130868 rows × 4 columns

```

1 trip_cr['year'].value_counts()
2 # Data contains order from only 2018

```

```

2018    130868
Name: year, dtype: int64

```

```

1 trip_cr['month'].value_counts()
2 # Data contains order from september and october month.

```

```

9     114979
10    15889
Name: month, dtype: int64

```

## ▼ Business Insights and Recommendations

```

1 # variables showing are highly right skew distribution
2 # mean of actual_time is higher than mean of osrm_time
3 # mean of actual_time is different than mean of segment_osrm_time
4 # mean of osrm_distance is similar as mean of segment_osrm_distance
5 # FTL Transportation used 69%
6 # Carting Transportation 31%
7 # There are 658 routes connecting sources to destinations
8 # most orders are coming from Bhindwandi_Mankoli_HB
9 # most orders are going to Mumbai
10 # the busiest route is Bhindwandi_Mankoli_HB to Mumbai
11 # company can use more transportation vehicles on this route
12 # most orders are coming from Maharashtra state
13 # most orders are going to Maharashtra state
14 # trip between cities Angamaly to Chalakudy saw the least avg time for completion
15 # trip between cities Hyderabad to Shamshabad saw the highest avg time for completion
16 # By visualizing plots, we observed that there is not much difference between pair variable
17 # All the appropriate insights and recommendations are mentioned with relevant code and plots

```