

▼ Problem Statement

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

In this case study we'll perform the following tasks

- Feature Engineering and feature creation
- Handling Missing values
- Clean, sanitize and perform EDA
- Perform hypothesis testing
- Prepare the data for further data processing pipeline and model building

```
1 !gdown 1lSJ5J1F-dN86PQTkQ_N0xWdIsbD0A1EX

  Downloading...
  From: https://drive.google.com/uc?id=1lSJ5J1F-dN86PQTkQ\_N0xWdIsbD0A1EX
  To: /content/delhivery_data.csv
  100% 55.6M/55.6M [00:00<00:00, 304MB/s]
```

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 pd.set_option('display.max_columns',None)
5 import seaborn as sns
6 from datetime import datetime as dt
7 from sklearn.preprocessing import OneHotEncoder
8 import scipy.stats as stats
9 import pylab
10 from sklearn import preprocessing
11 import warnings
12 warnings.filterwarnings("ignore")
13
```

```
1 !ls

delhivery_data.csv  sample_data
```

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

1 df.head(5)
```

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

```
1 df.shape

(144867, 24)
```

```
1 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

```

```
1 df.describe(include='all')
```

	data	trip_creation_time	route_schedule_uuid	route_type	trip_uuid	source_center	
count	144867	144867	144867	144867	144867	144867	
unique	2	14817	1504	2	14817	1508	
top	training	2018-09-28 05:23:15.359220	thanos::sroute:4029a8a2- 6c74-4b7e-a6d8-f9e069f...	FTL	153811219535896559	IND000000ACB	Gurga...
freq	104858	101	1812	99660	101	23347	
mean	NaN	NaN	NaN	NaN	NaN	NaN	NaN
std	NaN	NaN	NaN	NaN	NaN	NaN	NaN
min	NaN	NaN	NaN	NaN	NaN	NaN	NaN
25%	NaN	NaN	NaN	NaN	NaN	NaN	NaN
50%	NaN	NaN	NaN	NaN	NaN	NaN	NaN
75%	NaN	NaN	NaN	NaN	NaN	NaN	NaN
max	NaN	NaN	NaN	NaN	NaN	NaN	NaN

▼ Handling Missing Values

```

1 df.fillna('nan',inplace=True)
2
3 source_ids = list(df[df['source_name']=='nan']['source_center'].unique())
4 destination_ids = list(df[df['destination_name']=='nan']['destination_center'].unique())
5 source_ids.extend(destination_ids)
6 missing_ids = list(set(source_ids))
7 missing_ids

['IND282002AAD',
 'IND505326AAB',
 'IND465333A1B',
 'IND221005A1A',
 'IND577116AAA',
 'IND126116AAA',
 'IND122015AAC',
 'IND841301AAC',
 'IND250002AAC',
 'IND342902A1B',
 'IND509103AAC',
 'IND852118A1B',
 'IND331022A1B',
 'IND331001A1C']

1 len(missing_ids)

```

```

1 mapper = {}
2 counter=0
3 def missing(source_name,source_center,destination_name,destination_center):
4     if ((source_center in missing_ids) and (source_name!='nan')):
5         mapper[source_center]=source_name
6     if ((destination_center in missing_ids) and (destination_center!='nan')):
7         mapper[destination_center]=destination_name
8
9 df.apply(lambda x: missing(x.source_name,x.source_center,x.destination_name,x.destination_center),axis=1)
10 print(len(mapper))

```

0

```

1 counter=0
2
3 for a in missing_ids:
4     if a not in mapper:
5         mapper[a]='location_'+str(counter)+ '(' +str(counter)+ ')'
6         counter+=1
7 mapper

```

```
{
'IND282002AAD': 'location_0(0)',
'IND505326AAB': 'location_1(1)',
'IND465333A1B': 'location_2(2)',
'IND221005A1A': 'location_3(3)',
'IND577116AAA': 'location_4(4)',
'IND126116AAA': 'location_5(5)',
'IND122015AAC': 'location_6(6)',
'IND841301AAC': 'location_7(7)',
'IND250002AAC': 'location_8(8)',
'IND342902A1B': 'location_9(9)',
'IND509103AAC': 'location_10(10)',
'IND852118A1B': 'location_11(11)',
'IND331022A1B': 'location_12(12)',
'IND331001A1C': 'location_13(13)'}

```

```

1 def handling_missing(source_name,source_center,destination_name,destination_center):
2     if source_center in mapper and source_name == 'nan':
3         source_name=mapper[source_center]
4     if destination_center in mapper and destination_name == 'nan':
5         destination_name = mapper[destination_center]
6     return source_name,destination_name
7

```

```

1 df[['source_name','destination_name']] = df.apply(lambda x: pd.Series((
2     handling_missing(x.source_name,x.source_center,x.destination_name,x.destination_center)[0],
3     handling_missing(x.source_name,x.source_center,x.destination_name,x.destination_center)[1])),
4 axis=1)

```

```

1 df.isnull().sum()/df.shape[0]
2

```

data	0.0
trip_creation_time	0.0
route_schedule_uuid	0.0
route_type	0.0
trip_uuid	0.0
source_center	0.0
source_name	0.0
destination_center	0.0
destination_name	0.0
od_start_time	0.0
od_end_time	0.0
start_scan_to_end_scan	0.0
is_cutoff	0.0
cutoff_factor	0.0
cutoff_timestamp	0.0
actual_distance_to_destination	0.0
actual_time	0.0
osrm_time	0.0
osrm_distance	0.0
factor	0.0
segment_actual_time	0.0
segment_osrm_time	0.0
segment_osrm_distance	0.0
segment_factor	0.0

dtype: float64

▼ Feature Engineering

```

1 def segregate_city_state(data):
2     if data=='nan':
3         return 'nan','nan'
4     try:
5         location,state = data.split('(')[0][:-1],data.split('(')[1][:-1]
6         city = location.split('_')[0]
7         return city,state
8     except :
9         print(data)
10
11 df[['source_city','source_state']] = df.apply(lambda x: pd.Series((segregate_city_state(x.source_name)[0],
12                                         segregate_city_state(x.source_name)[1])),
13                                         axis=1)
14
15
16
17 df[['destination_city','destination_state']] = df.apply(lambda x: pd.Series((segregate_city_state(x.destination_name)[0]
18                                         ,segregate_city_state(x.destination_name)[1])),
19                                         axis=1)

```

```

1 def segregate_delivery_time(data):
2     year,month,date = data.split('-')[0],data.split('-')[1],data.split('-')[2][:3]
3     return year,month,date
4
5 df[['trip_creation_year','trip_creation_month','trip_creation_date']] = df.apply(lambda x: pd.Series((
6     segregate_delivery_time(x.trip_creation_time)[0],
7     segregate_delivery_time(x.trip_creation_time)[1],
8     segregate_delivery_time(x.trip_creation_time)[2])),
9     axis=1)

```

```
1 df['od_total_time'] = (pd.to_datetime(df['od_end_time'])-pd.to_datetime(df['od_start_time'])).astype('timedelta64[m]')

```

```
1 df.head(10)
```

	data	trip_creation_time	route_schedule_uuid	route_type	trip_uuid	source_center	sc
0	training	2018-09-20 02:35:36.476840	thanos::sroute:eb7bfc78- b351-4c0e-a951- fa3d5c3...	Carting	153741093647649320	trip- IND388121AAA	Anand_VI
1	training	2018-09-20 02:35:36.476840	thanos::sroute:eb7bfc78- b351-4c0e-a951- fa3d5c3...	Carting	153741093647649320	trip- IND388121AAA	Anand_VI
2	training	2018-09-20 02:35:36.476840	thanos::sroute:eb7bfc78- b351-4c0e-a951- fa3d5c3...	Carting	153741093647649320	trip- IND388121AAA	Anand_VI
3	training	2018-09-20 02:35:36.476840	thanos::sroute:eb7bfc78- b351-4c0e-a951- fa3d5c3...	Carting	153741093647649320	trip- IND388121AAA	Anand_VI
4	training	2018-09-20 02:35:36.476840	thanos::sroute:eb7bfc78- b351-4c0e-a951- fa3d5c3...	Carting	153741093647649320	trip- IND388121AAA	Anand_VI
5	training	2018-09-20 02:35:36.476840	thanos::sroute:eb7bfc78- b351-4c0e-a951- fa3d5c3...	Carting	153741093647649320	trip- IND388620AAB	Khamhat_M
6	training	2018-09-20 02:35:36.476840	thanos::sroute:eb7bfc78- b351-4c0e-a951- fa3d5c3...	Carting	153741093647649320	trip- IND388620AAB	Khamhat_M
7	training	2018-09-20 02:35:36.476840	thanos::sroute:eb7bfc78- b351-4c0e-a951- fa3d5c3...	Carting	153741093647649320	trip- IND388620AAB	Khamhat_M
8	training	2018-09-20 02:35:36.476840	thanos::sroute:eb7bfc78- b351-4c0e-a951- fa3d5c3...	Carting	153741093647649320	trip- IND388620AAB	Khamhat_M
9	training	2018-09-20 02:35:36.476840	thanos::sroute:eb7bfc78- b351-4c0e-a951- fa3d5c3...	Carting	153741093647649320	trip- IND388620AAB	Khamhat_M

▼ Data Cleaning

```

1 # removing columns which are not required
2 remove_columns = ['data','od_start_time','od_end_time','route_schedule_uuid','source_name',

```

```

3             'destination_name','is_cutoff','cutoff_factor','cutoff_timestamp','factor',
4             'trip_creation_time','segment_factor']
5 df.drop(remove_columns,axis=1,inplace=True)
6 df.shape
(144867, 20)

```

```

1 column_order = ['trip_uuid','route_type','source_center','destination_center','source_state',
2                  'destination_state','source_city','destination_city','trip_creation_year','trip_creation_month',
3                  'trip_creation_date','start_scan_to_end_scan','od_total_time','actual_distance_to_destination',
4                  'actual_time','osrm_time','osrm_distance','segment_actual_time','segment_osrm_time',
5                  'segment_osrm_distance']
6 df = df[column_order]

```

```
1 df.head()
```

	trip_uuid	route_type	source_center	destination_center	source_state	destination_state	source_c	
0	153741093647649320	trip-	Carting	IND388121AAA	IND388620AAB	Gujarat	Gujarat	An
1	153741093647649320	trip-	Carting	IND388121AAA	IND388620AAB	Gujarat	Gujarat	An
2	153741093647649320	trip-	Carting	IND388121AAA	IND388620AAB	Gujarat	Gujarat	An
3	153741093647649320	trip-	Carting	IND388121AAA	IND388620AAB	Gujarat	Gujarat	An
4	153741093647649320	trip-	Carting	IND388121AAA	IND388620AAB	Gujarat	Gujarat	An

▼ Merging Columns

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

```

1 data['segment_actual_time_cum'] = data.groupby(['trip_uuid','source_center',
2                                                 'destination_center'])['segment_actual_time'].cumsum()
3 data['segment_osrm_time_cum'] = data.groupby(['trip_uuid','source_center',
4                                                 'destination_center'])['segment_osrm_time'].cumsum()
5 data['segment_osrm_distance_cum'] = data.groupby(['trip_uuid','source_center',
6                                                 'destination_center'])['segment_osrm_distance'].cumsum()

```

```
1 data.head()
```

	trip_uuid	route_type	source_center	destination_center	source_state	destination_state	source_c	
0	153741093647649320	trip-	Carting	IND388121AAA	IND388620AAB	Gujarat	Gujarat	An
1	153741093647649320	trip-	Carting	IND388121AAA	IND388620AAB	Gujarat	Gujarat	An
2	153741093647649320	trip-	Carting	IND388121AAA	IND388620AAB	Gujarat	Gujarat	An
3	153741093647649320	trip-	Carting	IND388121AAA	IND388620AAB	Gujarat	Gujarat	An
4	153741093647649320	trip-	Carting	IND388121AAA	IND388620AAB	Gujarat	Gujarat	An

```

1 cumulative_columns = ['route_type','trip_creation_year','trip_creation_month','trip_creation_date',
2                       'od_total_time','start_scan_to_end_scan','actual_distance_to_destination',
3                       'actual_time','osrm_time','osrm_distance','source_city','destination_city',
4                       'source_state','destination_state','segment_actual_time_cum','segment_osrm_time_cum',
5                       'segment_osrm_distance_cum']
6 agg_on_trip_source_dest = data.groupby(['trip_uuid','source_center','destination_center']).last()[cumulative_columns]

```

```
1 agg_on_trip_source_dest.head(10)
```

			route_type	trip_creation_year	trip_creation_month	trip_creation_date	od_
	trip_uuid	source_center	destination_center				
trip- 153671041653548748	IND209304AAA	IND000000ACB	FTL	2018	09	12	
	IND462022AAA	IND209304AAA	FTL	2018	09	12	
trip- 153671042288605164	IND561203AAB	IND562101AAA	Carting	2018	09	12	
	IND572101AAA	IND561203AAB	Carting	2018	09	12	
trip- 153671043369099517	IND000000ACB	IND160002AAC	FTL	2018	09	12	
	IND562132AAA	IND000000ACB	FTL	2018	09	12	
trip- 153671046011330457	IND400072AAB	IND401104AAA	Carting	2018	09	12	
trip- 153671046011330457	IND583101AAA	IND583201AAA	FTL	2018	09	12	

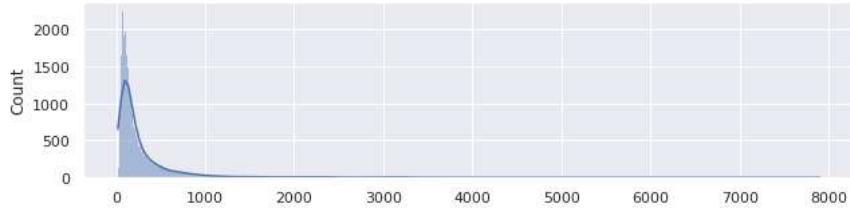
▼ Analysis

▼ Distribution of continues variables

```

1 sns.set()
2 fig,axiss= plt.subplots(5,1,figsize=(10,15))
3 sns.histplot(data = agg_on_trip_source_dest,x='start_scan_to_end_scan',kde=True,ax=axiss[0])
4 sns.histplot(data=agg_on_trip_source_dest,x='actual_distance_to_destination',kde=True,ax=axiss[1])
5 sns.histplot(data=agg_on_trip_source_dest,x='actual_time',kde=True,ax=axiss[2])
6 sns.histplot(data=agg_on_trip_source_dest,x='osrm_time',kde=True,ax=axiss[3])
7 sns.histplot(data=agg_on_trip_source_dest,x='osrm_distance',kde=True,ax=axiss[4])
8
9 plt.subplots_adjust(left=0.1,
10                     bottom=0.1,
11                     right=0.9,
12                     top=0.9,
13                     wspace=0.4,
14                     hspace=0.4)
15 plt.show()

```



▼ Outlier detection and Handling

```
1 agg_on_trip_source_dest.columns  
  
Index(['route_type', 'trip_creation_year', 'trip_creation_month',  
       'trip_creation_date', 'od_total_time', 'start_scan_to_end_scan',  
       'actual_distance_to_destination', 'actual_time', 'osrm_time',  
       'osrm_distance', 'source_city', 'destination_city', 'source_state',  
       'destination_state', 'segment_actual_time_cum', 'segment_osrm_time_cum',  
       'segment_osrm_distance_cum'],  
      dtype='object')  
  
1 sns.set()  
2 fig,axiss= plt.subplots(5,1,figsize=(10,15))  
3 sns.boxplot(data = agg_on_trip_source_dest,x='start_scan_to_end_scan',ax=ax[0])  
4 sns.boxplot(data=agg_on_trip_source_dest,x='actual_distance_to_destination',ax=ax[1])  
5 sns.boxplot(data=agg_on_trip_source_dest,x='actual_time',ax=ax[2])  
6 sns.boxplot(data=agg_on_trip_source_dest,x='osrm_time',ax=ax[3])  
7 sns.boxplot(data=agg_on_trip_source_dest,x='osrm_distance',ax=ax[4])  
8  
9 plt.subplots_adjust(left=0.1,  
10                     bottom=0.1,  
11                     right=0.9,  
12                     top=0.9,  
13                     wspace=0.4,  
14                     hspace=0.4)  
15 plt.show()
```



```

1 def detect_outlier(df, col_name):
2     q1 = df[col_name].quantile(0.25)
3     q3 = df[col_name].quantile(0.75)
4     iqr = q3-q1 #Interquartile range
5     fence_low = q1-1.5*iqr
6     fence_high = q3+1.5*iqr
7     df_out = df.loc[(df[col_name] < fence_low) | (df[col_name] > fence_high)]
8     return df_out
9
10
11 outlier_start_scan_to_end_scan = detect_outlier(agg_on_trip_source_dest,'start_scan_to_end_scan')
12 outlier_start_scan_to_end_scan

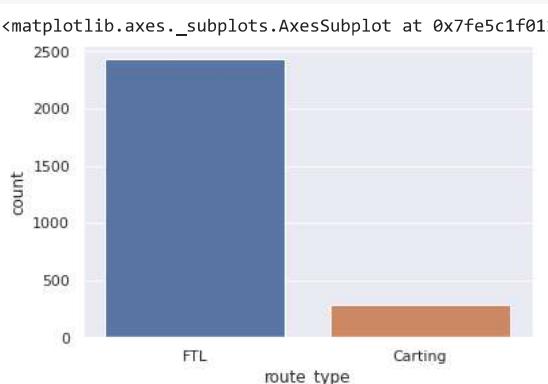
```

		route_type	trip_creation_year	trip_creation_month	tri
	trip_uuid	source_center	destination_center		
153671041653548748	trip-	IND209304AAA	IND000000ACB	FTL	2018 09
		IND462022AAA	IND209304AAA	FTL	2018 09
153671043369099517	trip-	IND000000ACB	IND160002AAC	FTL	2018 09
		IND562132AAA	IND000000ACB	FTL	2018 09
153671186247781647	trip-	IND209801AAA	IND209304AAA	Carting	2018 09

153860879439383883	trip-	IND562132AAA	IND000000ACB	FTL	2018 10
153860880135634048	trip-	IND425405AAA	IND424006AAA	FTL	2018 10
153860922975807074	trip-	IND508207AAB	IND507002AAA	FTL	2018 10
153861014185597051	trip-	IND462022AAA	IND209304AAA	FTL	2018 10
153861059679001096	trip-	IND208012AAA	IND209304AAA	Carting	2018 10

2720 rows × 17 columns

```
1 sns.countplot(x='route_type',data = outlier_start_scan_to_end_scan)
```



```

1 outlier_actual_distance_to_destination = detect_outlier(agg_on_trip_source_dest,'actual_distance_to_destination')
2 outlier_actual_distance_to_destination

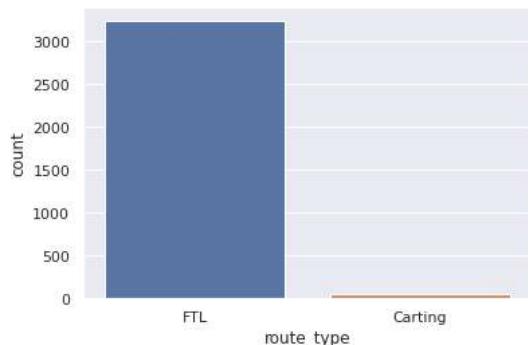
```

			route_type	trip_creation_year	trip_creation_month	tri
	trip_uuid	source_center	destination_center			
trip- 153671041653548748	IND209304AAA	IND000000ACB	FTL	2018	09	
	IND462022AAA	IND209304AAA	FTL	2018	09	
trip- 153671043369099517	IND000000ACB	IND160002AAC	FTL	2018	09	
	IND562132AAA	IND000000ACB	FTL	2018	09	
trip- 153671121411074590	IND501359AAE	IND515004AAA	FTL	2018	09	

trip- 153860840187622919	IND444005AAB	IND421302AAG	FTL	2018	10	
trip-	IND000000ACB	IND160002AAC	FTL	2018	10	

```
1 sns.countplot(x='route_type',data = outlier_actual_distance_to_destination)
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7fe5c88f13d0>
```



```
1 outlier_actual_time = detect_outlier(agg_on_trip_source_dest,'actual_time')
2 outlier_actual_time
```

			route_type	trip_creation_year	trip_creation_month	tri
	trip_uuid	source_center	destination_center			
trip- 153671041653548748	IND209304AAA	IND000000ACB	FTL	2018	09	
	IND462022AAA	IND209304AAA	FTL	2018	09	
trip- 153671043369099517	IND000000ACB	IND160002AAC	FTL	2018	09	
	IND562132AAA	IND000000ACB	FTL	2018	09	
trip- 153671121411074590	IND501359AAE	IND515004AAA	FTL	2018	09	

trip- 153860880135634048	IND425405AAA	IND424006AAA	FTL	2018	10	
trip- 153861007249500192	IND847404AAB	IND842001AAA	FTL	2018	10	
trip- 153861014185597051	IND206001AAA	IND000000ACB	FTL	2018	10	
	IND462022AAA	IND209304AAA	FTL	2018	10	
trip- 153861059679001096	IND208012AAA	IND209304AAA	Carting	2018	10	

3152 rows × 17 columns

```
1 sns.countplot(x='route_type',data = outlier_actual_time)
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7fe5cb37b430>
```

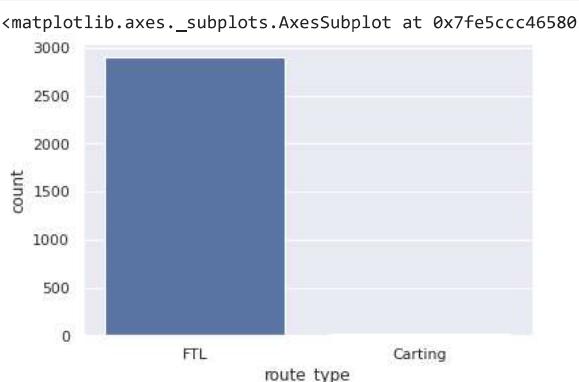


```
1 outlier_osrm_time = detect_outlier(agg_on_trip_source_dest,'osrm_time')
2 outlier_osrm_time
```

```
          route_type trip_creation_year trip_creation_month tri
trip_uuid source_center destination_center
153671041653548748 IND209304AAA IND000000ACB FTL 2018 09
                           IND462022AAA IND209304AAA FTL 2018 09
153671043369099517 IND000000ACB IND160002AAC FTL 2018 09
                           IND562132AAA IND000000ACB FTL 2018 09
153671121411074590 IND501359AAE IND515004AAA FTL 2018 09
...
153860840187622919 IND444005AAB IND421302AAG FTL 2018 10
153860879439383883 IND000000ACB IND160002AAC FTL 2018 10
                           IND562132AAA IND000000ACB FTL 2018 10
153861014185597051 IND206001AAA IND000000ACB FTL 2018 10
                           IND462022AAA IND209304AAA FTL 2018 10
```

2919 rows × 17 columns

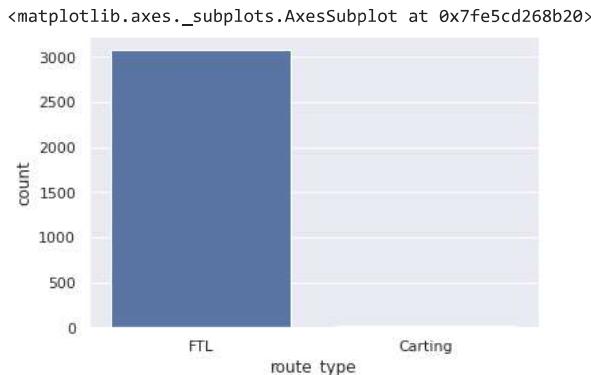
```
1 sns.countplot(x='route_type',data = outlier_osrm_time)
```



```
1 outlier_osrm_distance = detect_outlier(agg_on_trip_source_dest,'osrm_distance')
2 outlier_osrm_distance
```

			route_type	trip_creation_year	trip_creation_month	tri
	trip_uuid	source_center	destination_center			
153671041653548748	trip- IND209304AAA IND462022AAA	IND000000ACB IND209304AAA	IND000000ACB IND160002AAC	FTL FTL	2018 2018	09 09
	trip-			FTL	2018	09

```
1 sns.countplot(x='route_type',data = outlier_osrm_distance)
```



Conclusion: Based on above result We can see almost all the outlier data points belong to FTL route type and FTL being long journeys will usually take long time. Since all the outliers are valid values, we are going to consider them in our analysis

```
1 busiest_city_corridors = agg_on_trip_source_dest.groupby(['source_city',
2                                         'destination_city']).size().reset_index().rename(columns={0:'count'})
3 busiest_city_corridors.sort_values(by='count',ascending=False).reset_index().head(10)
```

	index	source_city	destination_city	count
0	272	Bengaluru	Bengaluru	565
1	192	Bangalore	Bengaluru	492
2	340	Bhiwandi	Mumbai	407
3	270	Bengaluru	Bangalore	356
4	958	Hyderabad	Hyderabad	316
5	1530	Mumbai	Mumbai	286
6	1528	Mumbai	Bhiwandi	282
7	579	Delhi	Gurgaon	248
8	481	Chennai	Chennai	246
9	815	Gurgaon	Delhi	237

Conclusion: Based on above table we can see the top busiest routes which includes bangalore, mumbai and hyderabad.

```
1 busiest_city_corridors_dist_time = agg_on_trip_source_dest.groupby(['source_city',
2                                         'destination_city'])['actual_time',
3                                         'actual_distance_to_destination'].agg('mean').reset_index()
4 pd.merge(busiest_city_corridors_dist_time,
5         busiest_city_corridors, left_index=True,right_index=True).drop(['source_city_y',
6         'destination_city_y'],axis=1).sort_values(by='count',
7         ascending=False).reset_index().rename({'source_city_x':'source_city',
8         'destination_city_x':'destination_city'})
```

index	source_city_x	destination_city_x	actual_time	actual_distance_to_destination	count
0	272	Bengaluru	Bengaluru	79.277876	29.075026 565
1	192	Bangalore	Bengaluru	77.674797	27.745791 492
2	340	Bhiwandi	Mumbai	80.120393	22.560310 407

```

1 busiest_state_corridors = agg_on_trip_source_dest.groupby(['source_state',
2                                         'destination_state']).size().reset_index().rename(columns={0:'count'})
3 busiest_state_corridors.sort_values(by='count',ascending=False).reset_index().head(10)

```

index	source_state	destination_state	count
0	120	Maharashtra	Maharashtra 3255
1	93	Karnataka	Karnataka 3158
2	153	Tamil Nadu	Tamil Nadu 2021
3	172	Uttar Pradesh	Uttar Pradesh 1526
4	163	Telangana	Telangana 1315
5	182	West Bengal	West Bengal 1296
6	50	Gujarat	Gujarat 1279
7	12	Andhra Pradesh	Andhra Pradesh 1139
8	148	Rajasthan	Rajasthan 1054
9	29	Bihar	Bihar 1011

```

1 busiest_state_corridors_dist_time = agg_on_trip_source_dest.groupby(['source_state',
2                                         'destination_state'])['actual_time',
3                                         'actual_distance_to_destination'].agg('mean').reset_index()
4
5 pd.merge(busiest_state_corridors,busiest_state_corridors_dist_time,
6         left_index=True,right_index=True).drop(['source_state_y','destination_state_y'],
7                                         axis=1).sort_values(by='count',
8                                         ascending=False).reset_index().rename(
9     {'source_city_x':'source_city','destination_city_x':'destination_city'}).head(10)

```

index	source_state_x	destination_state_x	count	actual_time	actual_distance_to_destination
0	120	Maharashtra	Maharashtra 3255	129.477419	48.333074
1	93	Karnataka	Karnataka 3158	92.926852	38.058320
2	153	Tamil Nadu	Tamil Nadu 2021	75.818407	34.846109
3	172	Uttar Pradesh	Uttar Pradesh 1526	136.305374	48.382655
4	163	Telangana	Telangana 1315	97.673004	42.062285
5	182	West Bengal	West Bengal 1296	136.176698	37.334884
6	50	Gujarat	Gujarat 1279	99.157936	48.589073
7	12	Andhra Pradesh	Andhra Pradesh 1139	97.304653	46.657201
8	148	Rajasthan	Rajasthan 1054	139.261860	64.128019
9	29	Bihar	Bihar 1011	163.715134	49.206346

Conclusion: Based on above table we can see the top busiest state routes which includes Maharashtra, Karnataka and Tamil Nadu

Handling Categorical values

```

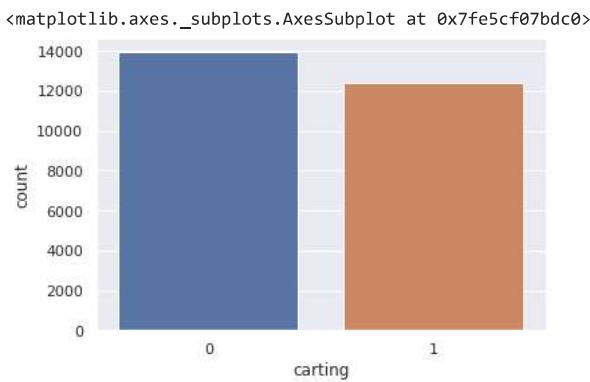
1 agg_on_trip_source_dest_encoder = pd.get_dummies(agg_on_trip_source_dest.route_type, prefix='route')
2 agg_on_trip_source_dest = agg_on_trip_source_dest.join(agg_on_trip_source_dest_encoder)
3 agg_on_trip_source_dest.rename({'route_Carting':'carting'},axis=1,inplace=True)
4 agg_on_trip_source_dest.drop('route_FTL',axis=1,inplace=True)
5 agg_on_trip_source_dest.drop('route_type',axis=1,inplace=True)
6 agg_on_trip_source_dest

```

			trip_creation_year	trip_creation_month	trip_creation_day
	trip_uuid	source_center	destination_center		
153671041653548748	trip-153671041653548748	IND209304AAA	IND000000ACB	2018	09
		IND462022AAA	IND209304AAA	2018	09
153671042288605164	trip-153671042288605164	IND561203AAB	IND562101AAA	2018	09
		IND572101AAA	IND561203AAB	2018	09
153671043369099517	trip-153671043369099517	IND000000ACB	IND160002AAC	2018	09

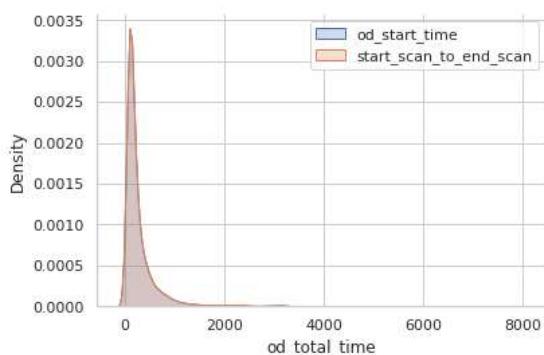
153861115439069069	trip-153861115439069069	IND628204AAA	IND627657AAA	2018	10
		IND628613AAA	IND627005AAA	2018	10
		IND628801AAA	IND628204AAA	2018	10

```
1 sns.countplot(data=agg_on_trip_source_dest,x='carting')
```



▼ 2 sample T test to check if od_total_time and start_scan_to_end_scan are same or not

```
1 sns.set_style('whitegrid')
2 sns.kdeplot(data=agg_on_trip_source_dest, x='od_total_time',fill=True,label='od_start_time')
3 sns.kdeplot(data=agg_on_trip_source_dest,x='start_scan_to_end_scan',fill=True,label='start_scan_to_end_scan')
4 sns.despine()
5 plt.legend()
6 plt.show()
7
```



Observation: From the above plot it is clear that the distribution is not normal for both the groups. Also we can see that od_start_time and start_scan_to_end_scan completely overlap each other.

Assumptions

- Both groups are independent
- Both groups are obtained through random sampling
- Data in each group is normally distributed
- variance of both the groups should be similar
- By visual analysis we can see that the data is not normally distributed for both groups(i.e od_total_time and start_scan_to_end_scan)

Hypothesis Formulation and test selection

- We'll have the following hypothesis
 - Null Hypothesis: The means of od_start_time and start_scan_to_end_scan are similar.
 - Alternate Hypothesis: The means of od_start_time and start_scan_to_end_scan are different.
 - We'll consider the significance value as 5% and perform a two tailed test
- Test Selection
 - We'll use 2 sample T-test since we need to compare mean of two independent group. The 2 sample t-test behaves similar to 2 sample z-test for large dataset(i.e. n>30)

▼ Checking Test assumptions

- We know that both groups are independent of each other since each both the times are calculated independently.
- We assume that both the groups are obtained from random sampling
- Data in each group is normally distributed - This assumption breaks as we have seen in the above plot. We need to apply a log transform to convert it to gaussian.
- Variance of both the groups must be similar. Below we'll see that the variance is somewhat similar for both groups

```

1 od_total_time_var = np.var(agg_on_trip_source_dest['od_total_time'])
2 start_scan_to_end_scan_variance = np.var(agg_on_trip_source_dest['start_scan_to_end_scan'])
3 print("variance of od_total_time is {} and variance of start_scan_to_end_scan is {}".format(od_total_time_var,
4                                                                                           start_scan_to_end_scan_variance))
5
variance of od_total_time is 194083.04130631435 and variance of start_scan_to_end_scan is 194083.04130631435

```

▼ Applying Log normal transformation to convert the data to gaussian

```

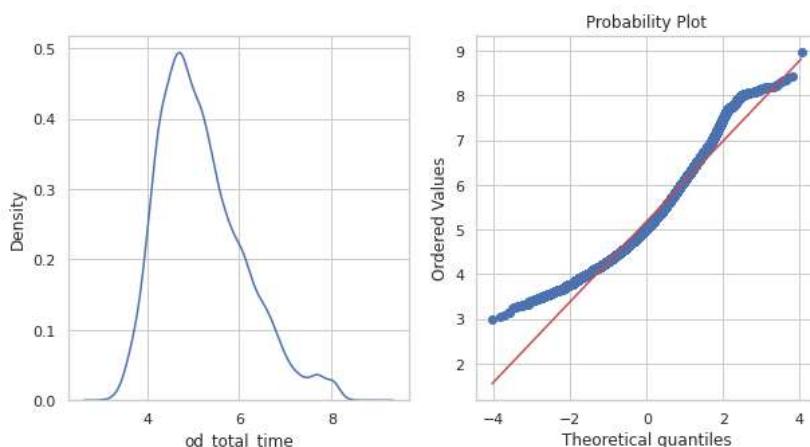
1 def normality(data):
2     plt.figure(figsize=(10,5))
3     plt.subplot(1,2,1)
4     sns.kdeplot(data)
5     plt.subplot(1,2,2)
6     stats.probplot(data,plot=pylab)
7     plt.show()

```

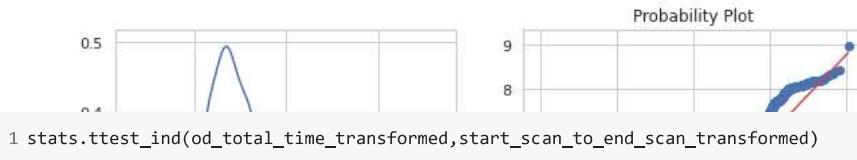
```

1 od_total_time_transformed = np.log(agg_on_trip_source_dest['od_total_time'])
2 start_scan_to_end_scan_transformed = np.log(agg_on_trip_source_dest['start_scan_to_end_scan'])
3 normality(od_total_time_transformed)

```



```
1 normality(start_scan_to_end_scan_transformed)
```



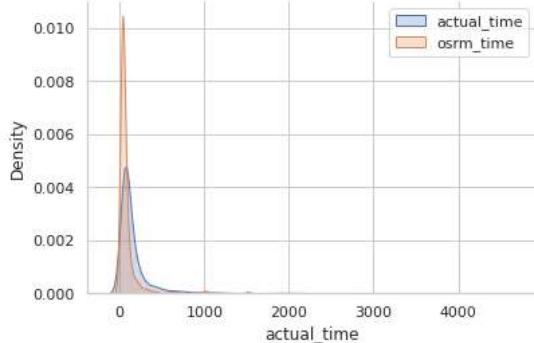
Conclusion: Based on above result p-value = 1 which is greater than our significance value alpha. So based on this We can say that we fail to reject the Null Hypothesis. It means that there is no difference between od_total_time and start_scan_to_end_scan.

2 sample T test to check if actual_time and osrm_time are same or not

```

1 sns.set_style('whitegrid')
2 sns.kdeplot(data=agg_on_trip_source_dest, x='actual_time', fill=True, label='actual_time')
3 sns.kdeplot(data=agg_on_trip_source_dest, x='osrm_time', fill=True, label='osrm_time')
4 sns.despine()
5 plt.legend()
6 plt.show()

```



Observation: From the above plot it is clear that the distribution is not normal for both the groups. Also we can see that actual_time and osrm_time don't completely overlap each other.

Assumptions

- Both groups are independent
- Both groups are obtained through random sampling
- Data in each group is normally distributed
- Variance of both the groups should be similar
- By visual analysis we can see that the data is not normally distributed for both groups(i.e actual_time and osrm_time)

Hypothesis Formulation and test selection

- We'll have the following hypothesis
 - Null Hypothesis: The means of actual_time and osrm_time are similar.
 - Alternate Hypothesis: The means of actual_time and osrm_time are different.
 - We'll consider the significance value as 5% and perform a two tailed test
- Test Selection
 - We'll use 2 sample T-test since we need to compare mean of two independent group. The 2 sample t-test behaves similar to 2 sample z-test for large dataset(i.e. n>30)

Checking Test assumptions

- We know that both groups are independent of each other since each both the times are calculated independently.
- We assume that both the groups are obtained from random sampling
- Data in each group is normally distributed - This assumption breaks as we have seen in the above plot. We need to apply a log transform to convert it to gaussian.
- Variance of both the groups must be similar. Below we'll see that the variance is somewhat similar for both groups

```

1 actual_time_var = np.var(agg_on_trip_source_dest['actual_time'])
2 osrm_time_variance = np.var(agg_on_trip_source_dest['osrm_time'])

```

```

3 print("variance of actual_time_var is {} and variance of osrm_time_variance is {}".format(actual_time_var,
4   variance of actual_time_var is 148106.70707614612 and variance of osrm_time_variance is 94291148400403046e))

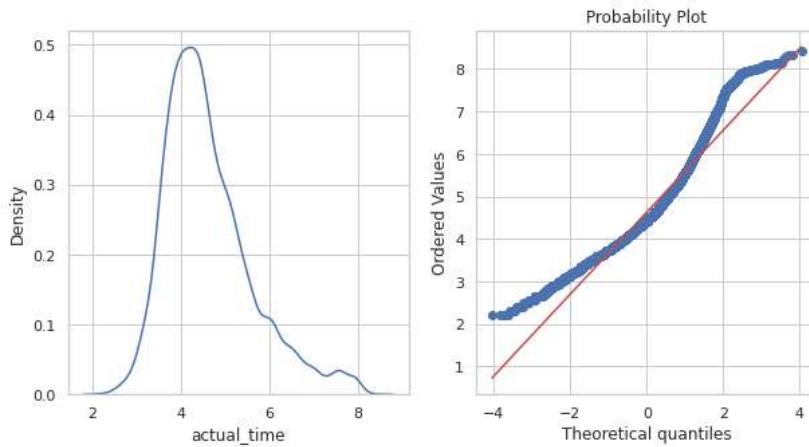
```

▼ Applying Log normal transformation to convert the data to gaussian

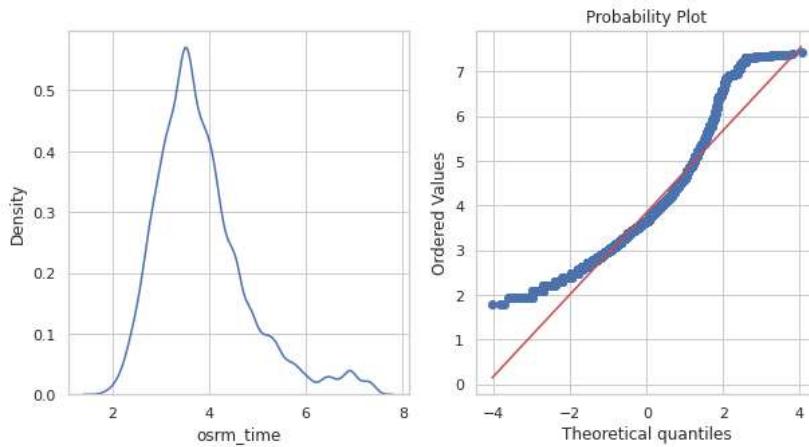
```

1 actual_time_transformed = np.log(agg_on_trip_source_dest['actual_time'])
2 osrm_time_transformed = np.log(agg_on_trip_source_dest['osrm_time'])
3 normality(actual_time_transformed)

```



```
1 normality(osrm_time_transformed)
```



```
1 stats.ttest_ind(actual_time_transformed,osrm_time_transformed)
```

```
Ttest_indResult(statistic=91.75651556695851, pvalue=0.0)
```

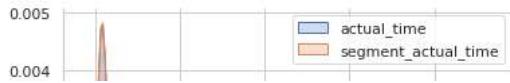
Conclusion: Based on above result p-value = 0 which is less than our significance value alpha. So based on this We can say that we reject the Null Hypothesis. It means that there is actual_time and osrm_time are different.

▼ 2 sample T test to check if actual_time and segment_actual_time are same or not

```

1 sns.set_style('whitegrid')
2 sns.kdeplot(data=agg_on_trip_source_dest, x='actual_time', fill=True, label='actual_time')
3 sns.kdeplot(data=agg_on_trip_source_dest,x='segment_actual_time_cum',fill=True,label='segment_actual_time')
4 plt.legend()
5 sns.despine()
6 plt.show()

```



Observation: From the above plot it is clear that the distribution is not normal for both the groups. Also we can see that actual_time and segment_actual_time don't completely overlap each other.

Assumptions

- Both groups are independent
- Both groups are obtained through random sampling
- Data in each group is normally distributed
- variance of both the groups should be similar
- By visual analysis we can see that the data is not normally distributed for both groups(i.e actual_time and segment_actual_time)

Hypothesis Formulation and test selection

- We'll have the following hypothesis
 - Null Hypothesis: The means of actual_time and segment_actual_time are similar.
 - Alternate Hypothesis: The means of actual_time and segment_actual_time are different.
 - We'll consider the significance value as 5% and perform a two tailed test
- Test Selection
 - We'll use 2 sample T-test since we need to compare mean of two independent group. The 2 sample t-test behaves similar to 2 sample z-test for large dataset(i.e. n>30)

Checking Test assumptions

- We know that both groups are independent of each other since each both the times are calculated independently.
- We assume that both the groups are obtained from random sampling
- Data in each group is normally distributed - This assumption breaks as we have seen in the above plot. We need to apply a log transform to convert it to gaussian.
- Variance of both the groups must be similar. Below we'll see that the variance is somewhat similar for both groups

```

1 actual_time_var = np.var(agg_on_trip_source_dest['actual_time'])
2 segment_actual_time_variance = np.var(agg_on_trip_source_dest['segment_actual_time_cum'])
3 print("variance of actual_time_var is {} and variance of segment_actual_time_variance is {}".format(actual_time_var,
4                                                               segment_actual_time_variance))

variance of actual_time_var is 148106.70707614612 and variance of segment_actual_time_variance is 145371.38348607288

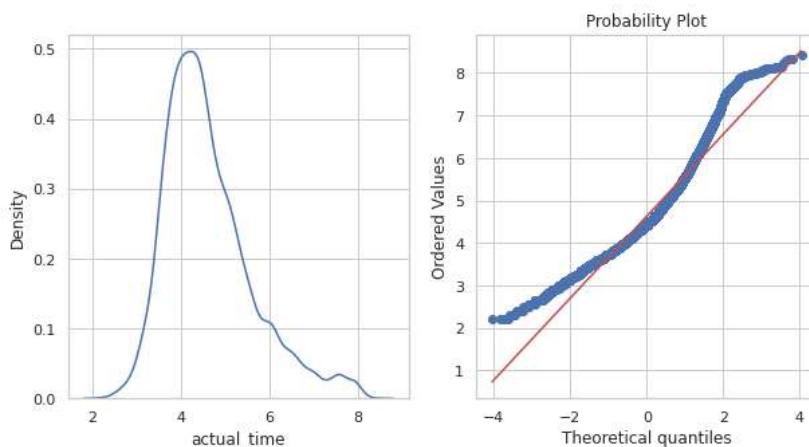
```

Applying Log normal transformation to convert the data to gaussian

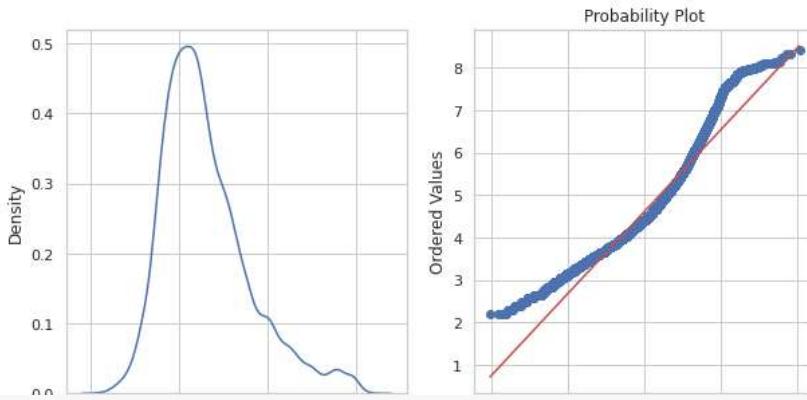
```

1 actual_time_transformed = np.log(agg_on_trip_source_dest['actual_time'])
2 segment_actual_time_transformed = np.log(agg_on_trip_source_dest['segment_actual_time_cum'])
3 normality(actual_time_transformed)

```



```
1 normality(segment_actual_time_transformed)
```

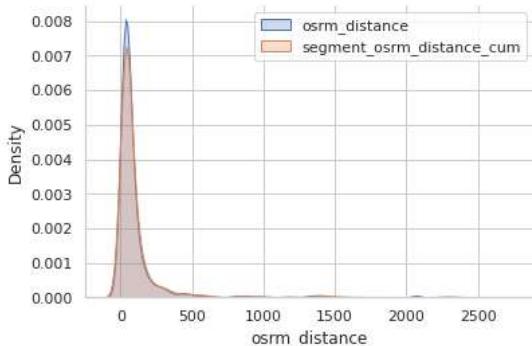


```
1 stats.ttest_ind(actual_time_transformed,segment_actual_time_transformed)
Ttest_indResult(statistic=1.1719659905322306, pvalue=0.24121600541669044)
```

Conclusion: Based on above result p-value = 0.24 which is greater than our significance value alpha. So based on this We can say that we fail to reject the Null Hypothesis. It means that there is actual_time and segment_actual_time are same.

▼ 2 sample T test to check if actual_time and segment_actual_time are same or not

```
1 sns.set_style('whitegrid')
2 sns.kdeplot(data=agg_on_trip_source_dest, x='osrm_distance',fill=True,label='osrm_distance')
3 sns.kdeplot(data=agg_on_trip_source_dest,x='segment_osrm_distance_cum',fill=True,label='segment_osrm_distance_cum')
4 plt.legend()
5 sns.despine()
6 plt.show()
```



Observation: From the above plot it is clear that the distribution is not normal for both the groups. Also we can see that osrm_distance and segment_osrm_distance_cum don't completely overlap each other.

Assumptions

- Both groups are independent
- Both groups are obtained through random sampling
- Data in each group is normally distributed
- variance of both the groups should be similar
- By visual analysis we can see that the data is not normally distributed for both groups(i.e osrm_distance and segment_osrm_distance_cum)

Hypothesis Formulation and test selection

- We'll have the following hypothesis
 - Null Hypothesis: The means of actual_time and segment_actual_time are similar.
 - Alternate Hypothesis: The means of osrm_distance and segment_osrm_distance_cum are different.
 - We'll consider the significance value as 5% and perform a two tailed test
- Test Selection
 - We'll use 2 sample T-test since we need to compare mean of two independent group. The 2 sample t-test behaves similar to 2 sample z-test for large dataset(i.e. n>30)

▼ Checking Test assumptions

- We know that both groups are independent of each other since each both the times are calculated independently.
- We assume that both the groups are obtained from random sampling
- Data in each group is normally distributed - This assumption breaks as we have seen in the above plot. We need to apply a log transform to convert it to gaussian.
- Variance of both the groups must be similar. Below we'll see that the variance is somewhat similar for both groups

```

1 osrm_distance_var = np.var(agg_on_trip_source_dest['osrm_distance'])
2 segment_osrm_distance_cum_var = np.var(agg_on_trip_source_dest['segment_osrm_distance_cum'])
3 print("variance of osrm_distance_var is {} and variance of segment_osrm_distance_cum_var is {}".format(
4     osrm_distance_var,segment_osrm_distance_cum_var))

variance of osrm_distance_var is 64398.681296934046 and variance of segment_osrm_distance_cum_var is 81754.32592611018

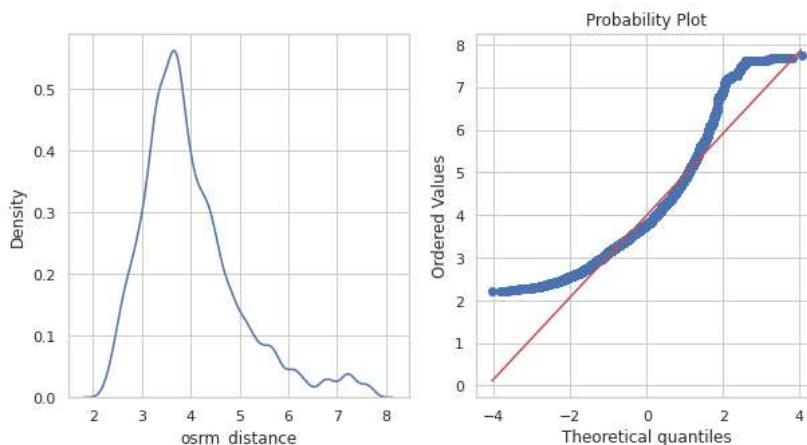
```

▼ Applying Log normal transformation to convert the data to gaussian

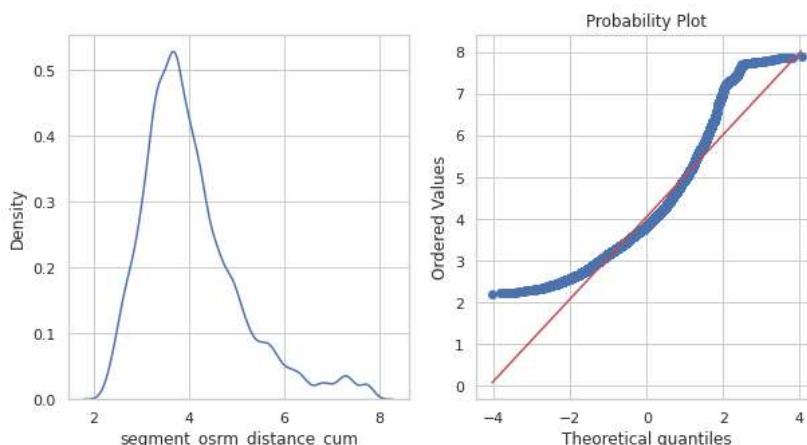
```

1 osrm_distance_transformed = np.log(agg_on_trip_source_dest['osrm_distance'])
2 segment_osrm_distance_cum_transformed = np.log(agg_on_trip_source_dest['segment_osrm_distance_cum'])
3 normality(osrm_distance_transformed)

```



```
1 normality(segment_osrm_distance_cum_transformed)
```



```
1 stats.ttest_ind(osrm_distance_transformed,segment_osrm_distance_cum_transformed)
```

```
Ttest_indResult(statistic=-6.016679309589272, pvalue=1.792010384210536e-09)
```

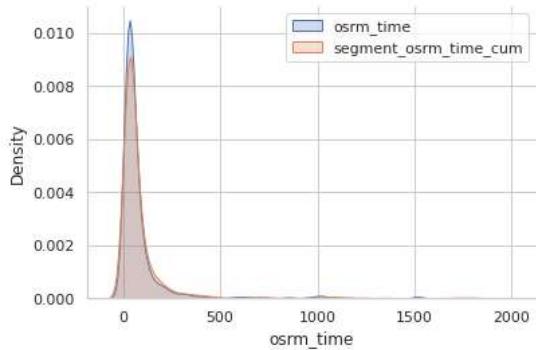
Conclusion: Based on above result p-value = 1.79^{-9} which is less than our significance value alpha. So based on this We can say that we reject the Null Hypothesis. It means that there is osrm_distance and segment_osrm_distance_cum are different.

▼ 2 sample T test to check if osrm_time_cum and segment_osrm_time_cum are same or not

```

1 sns.set_style('whitegrid')
2 sns.kdeplot(data=agg_on_trip_source_dest, x='osrm_time', fill=True, label='osrm_time')
3 sns.kdeplot(data=agg_on_trip_source_dest, x='segment_osrm_time_cum', fill=True, label='segment_osrm_time_cum')
4 plt.legend()
5 sns.despine()
6 plt.show()

```



Observation: From the above plot it is clear that the distribution is not normal for both the groups. Also we can see that osrm_time and segment_osrm_time_cum don't completely overlap each other.

Assumptions

- Both groups are independent
- Both groups are obtained through random sampling
- Data in each group is normally distributed
- variance of both the groups should be similar
- By visual analysis we can see that the data is not normally distributed for both groups(i.e osrm_time and segment_osrm_time_cum)

Hypothesis Formulation and test selection

- We'll have the following hypothesis
 - Null Hypothesis: The means of actual_time and segment_osrm_time_cum are similar.
 - Alternate Hypothesis: The means of osrm_time and segment_osrm_distance_cum are different.
 - We'll consider the significance value as 5% and perform a two tailed test
- Test Selection
 - We'll use 2 sample T-test since we need to compare mean of two independent group. The 2 sample t-test behaves similar to 2 sample z-test for large dataset(i.e. n>30)

Checking Test assumptions

- We know that both groups are independent of each other since each both the times are calculated independently.
- We assume that both the groups are obtained from random sampling
- Data in each group is normally distributed - This assumption breaks as we have seen in the above plot. We need to apply a log transform to convert it to gaussian.
- Variance of both the groups must be similar. Below we'll see that the variance is somewhat similar for both groups

```

1 osrm_time_var = np.var(agg_on_trip_source_dest['osrm_time'])
2 segment_osrm_time_var = np.var(agg_on_trip_source_dest['segment_osrm_time_cum'])
3 print("variance of osrm_time_var is {} and variance of segment_osrm_time_cum_var is {}".format(osrm_time_var,
4                                                                                           segment_osrm_time_var))

```

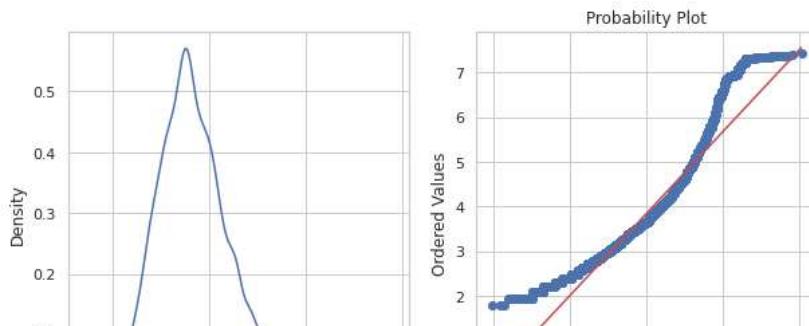
variance of osrm_time_var is 34253.46400403046 and variance of segment_osrm_time_cum_var is 46503.567775992764

Applying Log normal transformation to convert the data to gaussian

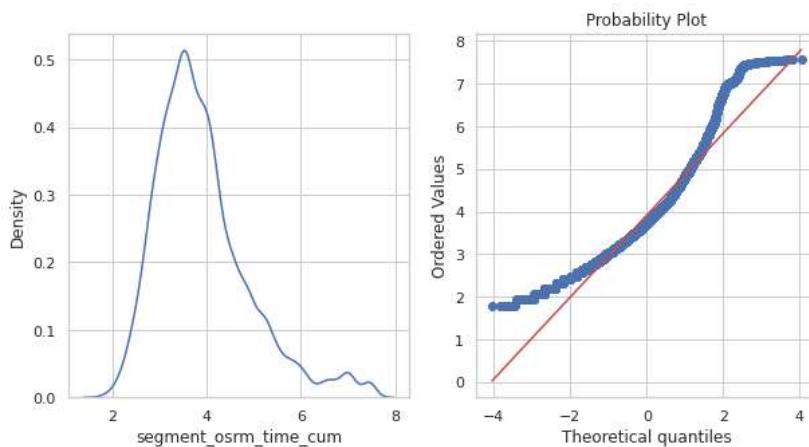
```

1 osrm_time_transformed = np.log(agg_on_trip_source_dest['osrm_time'])
2 segment_osrm_time_cum_transformed = np.log(agg_on_trip_source_dest['segment_osrm_time_cum'])
3 normality(osrm_time_transformed)

```



```
1 normality(segment_osrm_time_cum_transformed)
```



```
1 stats.ttest_ind(osrm_time_transformed,segment_osrm_time_cum_transformed)
```

```
Ttest_indResult(statistic=-7.264652594570266, pvalue=3.7915395546872833e-13)
```

Conclusion: Based on above result p-value = 3.3×10^{-13} which is less than our significance value alpha. So based on this We can say that we reject the Null Hypothesis. It means that there is osrm_time and segment_osrm_time_cum are different.

▼ Normalization

```
1 normalized_agg_on_trip_source_dest = agg_on_trip_source_dest.copy()
2 columns_to_normalize = ['od_total_time','start_scan_to_end_scan','actual_distance_to_destination',
3                         'actual_time','osrm_time','osrm_distance','segment_actual_time_cum',
4                         'segment_osrm_time_cum','segment_osrm_distance_cum']
5 normalized_agg_on_trip_source_dest[columns_to_normalize] = preprocessing.normalize(
6                                         agg_on_trip_source_dest[columns_to_normalize])
7 normalized_agg_on_trip_source_dest
```

			trip_creation_year	trip_creation_month	trip_creation_
	trip_uuid	source_center	destination_center		
trip- 153671041653548748	IND209304AAA	IND000000ACB	2018	09	
	IND462022AAA	IND209304AAA	2018	09	
trip- 153671042288605164	IND561203AAB	IND562101AAA	2018	09	
	IND572101AAA	IND561203AAB	2018	09	
trip- 153671043369099517	IND000000ACB	IND160002AAC	2018	09	

trip- 153861115439069069	IND628204AAA	IND627657AAA	2018	10	
	IND628613AAA	IND627005AAA	2018	10	
	IND628801AAA	IND628204AAA	2018	10	
trip- 153861118270144424	IND583119AAA	IND583101AAA	2018	10	
	IND583201AAA	IND583119AAA	2018	10	

26368 rows × 17 columns

▼ Recommendations

Based on above Analysis, following recommendations can be made:

- Most of the trips are taking place in tier 1 cities like bangalore, mumbai and hyderabad and are mostly inter city. There are not many interstate deliveries taking place. Team can focus on inter state deliveries since it can generate more revenue due to long distance parcels.
- The average delivery time in west bengal is more in comparison to other states. Team can focus on reducing this delivery time.
- osrm_time and segment_osrm_time are different. Team needs to make sure this difference is reduced, so that better delivery time prediction can be made and it becomes convenient for the customer to expect an accurate delivery time.
- The actual time taken to deliver and the one predicted by osrm are not same. Team needs to improve this prediction modes so as to give more accurate predictions.
- The osrm distance and actual distance covered are also not same i.e. maybe the delivery person is not following the predefined route which may lead to late deliveries or the osrm devices is not properly predicting the route based on distance, traffic and other factors. Team needs to look into it.