

```

import pandas as pd
import numpy as np
import seaborn as sns
from matplotlib import pyplot as plt
from sklearn import preprocessing

import warnings
warnings.filterwarnings('ignore')

df = pd.read_csv('delhi_delivery_data.csv')

df.shape

(144897, 24)

pd.set_option('display.max_columns', None)
df.head(10)

data
trip_creation_time route_schedule_uid route_type trip_uid source_center source_name destination_center destination_name od_start_time od_end_time start_scan_to_end_scan is_cutoff cutoff_factor cutoff_timestamp actual_distance
0 0 2018-09-20 02:35:36.476460 Phanos-route-e776c78-0351-40e-8b51-fa365c3... CarWng 153741003647649320 INDB88121AAA Anand_VUNagar_DC (Gurgaon) INDB88120AAB Khamhat_MotvDPP_D (Gurgaon) 2018-09-20 02:31:32.418600 04:47:45.236797 86.0 True 9 2018-09-20 04:27:55
1 1 2018-09-20 02:35:36.476460 Phanos-route-e776c78-0351-40e-8b51-fa365c3... CarWng 153741003647649320 INDB88121AAA Anand_VUNagar_DC (Gurgaon) INDB88120AAB Khamhat_MotvDPP_D (Gurgaon) 2018-09-20 02:31:32.418600 04:47:45.236797 86.0 True 18 2018-09-20 04:17:55
2 2 2018-09-20 02:35:36.476460 Phanos-route-e776c78-0351-40e-8b51-fa365c3... CarWng 153741003647649320 INDB88121AAA Anand_VUNagar_DC (Gurgaon) INDB88120AAB Khamhat_MotvDPP_D (Gurgaon) 2018-09-20 02:31:32.418600 04:47:45.236797 86.0 True 27 2018-09-20 04:01:19.555586
3 3 2018-09-20 02:35:36.476460 Phanos-route-e776c78-0351-40e-8b51-fa365c3... CarWng 153741003647649320 INDB88121AAA Anand_VUNagar_DC (Gurgaon) INDB88120AAB Khamhat_MotvDPP_D (Gurgaon) 2018-09-20 02:31:32.418600 04:47:45.236797 86.0 True 36 2018-09-20 03:59:57
4 4 2018-09-20 02:35:36.476460 Phanos-route-e776c78-0351-40e-8b51-fa365c3... CarWng 153741003647649320 INDB88121AAA Anand_VUNagar_DC (Gurgaon) INDB88120AAB Khamhat_MotvDPP_D (Gurgaon) 2018-09-20 02:31:32.418600 04:47:45.236797 86.0 False 39 2018-09-20 03:55:55
5 5 2018-09-20 02:35:36.476460 Phanos-route-e776c78-0351-40e-8b51-fa365c3... CarWng 153741003647649320 INDB88120AAB Khamhat_MotvDPP_D (Gurgaon) INDB88120AAB Anand_Vaghasi_IP (Gurgaon) 2018-09-20 04:47:45.236797 06:36:55.627764 109.0 True 9 2018-09-20 04:17:55
6 6 2018-09-20 02:35:36.476460 Phanos-route-e776c78-0351-40e-8b51-fa365c3... CarWng 153741003647649320 INDB88120AAB Khamhat_MotvDPP_D (Gurgaon) INDB88120AAB Anand_Vaghasi_IP (Gurgaon) 2018-09-20 04:47:45.236797 06:36:55.627764 109.0 True 18 2018-09-20 04:17:55
7 7 2018-09-20 02:35:36.476460 Phanos-route-e776c78-0351-40e-8b51-fa365c3... CarWng 153741003647649320 INDB88120AAB Khamhat_MotvDPP_D (Gurgaon) INDB88120AAB Anand_Vaghasi_IP (Gurgaon) 2018-09-20 04:47:45.236797 06:36:55.627764 109.0 True 27 2018-09-20 04:17:55
8 8 2018-09-20 02:35:36.476460 Phanos-route-e776c78-0351-40e-8b51-fa365c3... CarWng 153741003647649320 INDB88120AAB Khamhat_MotvDPP_D (Gurgaon) INDB88120AAB Anand_Vaghasi_IP (Gurgaon) 2018-09-20 04:47:45.236797 06:36:55.627764 109.0 True 36 2018-09-20 04:17:55
9 9 2018-09-20 02:35:36.476460 Phanos-route-e776c78-0351-40e-8b51-fa365c3... CarWng 153741003647649320 INDB88120AAB Khamhat_MotvDPP_D (Gurgaon) INDB88120AAB Anand_Vaghasi_IP (Gurgaon) 2018-09-20 04:47:45.236797 06:36:55.627764 109.0 False 43 2018-09-20 04:49:20

# Missing value detection
df.columns[df.isna().any()].tolist()
['source_name', 'destination_name']

# There would be a mapping with Delivery for source center to source name and destination center to destination name. That can be used to fill in the missing values. We do not have that mapping.
# INDB77116AAA appears as both source center and destination center and has no source or destination name corresponding to it. Using most frequent values also will lead to false data. Hence it is best to either delete such rows or map them to Other (India)
# To conclude given just Source Center we CANNOT DEDUCE correct and accurate Source Name from it

df = df.dropna(axis=1)
144316 rows x 24 new shape of data

df = df.reset_index(drop = True)

df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 144316 entries, 0 to 144315
Data columns (total 24 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   data        144316 non-null   object
 1   trip_creation_time  144316 non-null   object
 2   route_schedule_uid  144316 non-null   object
 3   route_type    144316 non-null   object
 4   trip_uid      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   is_cutoff      144316 non-null   bool
13   cutoff_factor  144316 non-null   int64
14   cutoff_timestamp  144316 non-null   object
15   actual_distance_to_destination  144316 non-null   float64
16   actual_time    144316 non-null   float64
17   osrm_time      144316 non-null   float64
18   Factor         144316 non-null   float64
19   segment_actual_time  144316 non-null   float64
20   segment_osrm_time  144316 non-null   float64
21   segment_osrm_distance  144316 non-null   float64
22   segment_factor  144316 non-null   float64
23   segment_factor  144316 non-null   float64
dtypes: bool(1), float64(19), int64(1), object(12)
memory usage: 25.5+ MB

# Subject to datetime conversion
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'])

# Categorical variables are data, route_type
df['data'] = df['data'].map({'training':1, 'test':0})
df['route_type'] = df['route_type'].map({'FTL':1, 'Carting':0})

df.head(20)

data
trip_creation_time route_schedule_uid route_type trip_uid source_center source_name destination_center destination_name od_start_time od_end_time start_scan_to_end_scan is_cutoff cutoff_factor cutoff_timestamp actual_distance
0 0 2018-09-20 02:35:36.476460 Phanos-route-e776c78-0351-40e-8b51-fa365c3... CarWng 153741003647649320 INDB88121AAA Anand_VUNagar_DC (Gurgaon) INDB88120AAB Khamhat_MotvDPP_D (Gurgaon) 2018-09-20 02:31:32.418600 04:47:45.236797 86.0 True 9 2018-09-20 04:27:55
1 1 2018-09-20 02:35:36.476460 Phanos-route-e776c78-0351-40e-8b51-fa365c3... CarWng 153741003647649320 INDB88121AAA Anand_VUNagar_DC (Gurgaon) INDB88120AAB Khamhat_MotvDPP_D (Gurgaon) 2018-09-20 02:31:32.418600 04:47:45.236797 86.0 True 18 2018-09-20 04:17:55
2 1 2018-09-20 02:35:36.476460 Phanos-route-e776c78-0351-40e-8b51-fa365c3... CarWng 153741003647649320 INDB88121AAA Anand_VUNagar_DC (Gurgaon) INDB88120AAB Khamhat_MotvDPP_D (Gurgaon) 2018-09-20 02:31:32.418600 04:47:45.236797 86.0 True 27 2018-09-20 04:01:19.555586
3 1 2018-09-20 02:35:36.476460 Phanos-route-e776c78-0351-40e-8b51-fa365c3... CarWng 153741003647649320 INDB88121AAA Anand_VUNagar_DC (Gurgaon) INDB88120AAB Khamhat_MotvDPP_D (Gurgaon) 2018-09-20 02:31:32.418600 04:47:45.236797 86.0 True 36 2018-09-20 03:59:57
4 1 2018-09-20 02:35:36.476460 Phanos-route-e776c78-0351-40e-8b51-fa365c3... CarWng 153741003647649320 INDB88121AAA Anand_VUNagar_DC (Gurgaon) INDB88120AAB Khamhat_MotvDPP_D (Gurgaon) 2018-09-20 02:31:32.418600 04:47:45.236797 86.0 False 39 2018-09-20 03:55:55
5 1 2018-09-20 02:35:36.
```