```
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib as mpl
import matplotlib.pyplot as plt
import scipy.stats as spy

import warnings
warnings.simplefilter('ignore')

df = pd.read_csv(r"https://d2beiqkhq929f0.cloudfront.net/public_assets/assets/000/001/551/original/delhivery_data.csv?1642751181")

df.head()
```

| | data | trip_creation_time | route_schedule_uuid | route_type | trip_uuid | sou |
|---|---|---|---|---|---|---|
| 0 | training | 2018-09-20 02:35:36.476840 | thanos::sroute:eb7bfc78-b351-4c0e-a951-fa3d5c3... | Carting | trip-1537410936476649320 | INI |
| 1 | training | 2018-09-20 02:35:36.476840 | thanos::sroute:eb7bfc78-b351-4c0e-a951-fa3d5c3... | Carting | trip-1537410936476649320 | INI |
| 2 | training | 2018-09-20 02:35:36.476840 | thanos::sroute:eb7bfc78-b351-4c0e-a951-fa3d5c3... | Carting | trip-1537410936476649320 | INI |
| 3 | training | 2018-09-20 02:35:36.476840 | thanos::sroute:eb7bfc78-b351-4c0e-a951-fa3d5c3... | Carting | trip-1537410936476649320 | INI |
| 4 | training | 2018-09-20 02:35:36.476840 | thanos::sroute:eb7bfc78-b351-4c0e-a951-fa3d5c3... | Carting | trip-1537410936476649320 | INI |

5 rows × 24 columns

```
df.shape
```

```
(144867, 24)
```

```
df.columns
```

```
Index(['data', 'trip_creation_time', 'route_schedule_uuid', 'route_type',
       'trip_uuid', 'source_center', 'source_name', 'destination_center',
       'destination_name', 'od_start_time', 'od_end_time',
       'start_scan_to_end_scan', 'is_cutoff', 'cutoff_factor',
       'cutoff_timestamp', 'actual_distance_to_destination', 'actual_time',
       'osrm_time', 'osrm_distance', 'factor', 'segment_actual_time',
       'segment_osrm_time', 'segment_osrm_distance', 'segment_factor'],
      dtype='object')
```

```
df.dtypes
```

```
data                            object
trip_creation_time              object
route_schedule_uuid             object
route_type                      object
trip_uuid                       object
source_center                   object
source_name                     object
destination_center              object
destination_name                object
od_start_time                   object
od_end_time                     object
start_scan_to_end_scan          float64
is_cutoff                         bool
cutoff_factor                    int64
cutoff_timestamp                object
actual_distance_to_destination  float64
actual_time                     float64
osrm_time                       float64
osrm_distance                   float64
factor                          float64
segment_actual_time             float64
```

```
segment_osrm_time              float64
segment_osrm_distance          float64
segment_factor                 float64
dtype: object
```

```
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
```

```
unknown_field = ['is_cutoff','cutoff_factor', 'cutoff_timestamp','factor','segment_factor']
df = df.drop(columns = unknown_field)
```

```
for i in df.columns:
  print(f"unique entries for column {i:<30} = {df[i].nunique()}")
```

```
unique entries for column data                           = 2
unique entries for column trip_creation_time             = 14817
unique entries for column route_schedule_uuid            = 1504
unique entries for column route_type                     = 2
unique entries for column trip_uuid                      = 14817
unique entries for column source_center                  = 1508
unique entries for column source_name                    = 1498
unique entries for column destination_center             = 1481
unique entries for column destination_name               = 1468
unique entries for column od_start_time                  = 26369
unique entries for column od_end_time                    = 26369
unique entries for column start_scan_to_end_scan         = 1915
unique entries for column actual_distance_to_destination = 144515
unique entries for column actual_time                    = 3182
unique entries for column osrm_time                      = 1531
unique entries for column osrm_distance                  = 138046
unique entries for column segment_actual_time            = 747
unique entries for column segment_osrm_time              = 214
unique entries for column segment_osrm_distance          = 113799
```

```
df['data'] = df['data'].astype('category')
df['route_type'] = df['route_type'].astype('category')
```

```
floating_columns = ['actual_distance_to_destination', 'actual_time', 'osrm_time', 'osrm_distance',
                    'segment_actual_time', 'segment_osrm_time', 'segment_osrm_distance']
for i in floating_columns:
    print(df[i].max())
```

```
1927.4477046975032
4532.0
1686.0
2326.1991000000003
3051.0
```

```
       1611.0
       2191.4037000000003


for i in floating_columns:
    df[i] = df[i].astype('float32')


datetime_columns = ['trip_creation_time', 'od_start_time', 'od_end_time']
for i in datetime_columns:
    df[i] = pd.to_datetime(df[i])


df.info()
```

```
       <class 'pandas.core.frame.DataFrame'>
       RangeIndex: 144867 entries, 0 to 144866
       Data columns (total 19 columns):
        #   Column                        Non-Null Count   Dtype
       ---  ------                        --------------   -----
        0   data                          144867 non-null  category
        1   trip_creation_time            144867 non-null  datetime64[ns]
        2   route_schedule_uuid           144867 non-null  object
        3   route_type                    144867 non-null  category
        4   trip_uuid                     144867 non-null  object
        5   source_center                 144867 non-null  object
        6   source_name                   144574 non-null  object
        7   destination_center            144867 non-null  object
        8   destination_name              144606 non-null  object
        9   od_start_time                 144867 non-null  datetime64[ns]
        10  od_end_time                   144867 non-null  datetime64[ns]
        11  start_scan_to_end_scan        144867 non-null  float64
        12  actual_distance_to_destination 144867 non-null float32
        13  actual_time                   144867 non-null  float32
        14  osrm_time                     144867 non-null  float32
        15  osrm_distance                 144867 non-null  float32
        16  segment_actual_time           144867 non-null  float32
        17  segment_osrm_time             144867 non-null  float32
        18  segment_osrm_distance         144867 non-null  float32
       dtypes: category(2), datetime64[ns](3), float32(7), float64(1), object(6)
       memory usage: 15.2+ MB


df['trip_creation_time'].min(), df['od_end_time'].max()
```

```
       (Timestamp('2018-09-12 00:00:16.535741'),
        Timestamp('2018-10-08 03:00:24.353479'))


np.any(df.isnull())
```

```
       True


df.isnull().sum()
```

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


missing_source_name = df.loc[df['source_name'].isnull(), 'source_center'].unique()
missing_source_name
```

```
       array(['IND342902A1B', 'IND577116AAA', 'IND282002AAD', 'IND465333A1B',
              'IND841301AAC', 'IND509103AAC', 'IND126116AAA', 'IND331022A1B',
              'IND505326AAB', 'IND852118A1B'], dtype=object)
```

```python
for i in missing_source_name:
    unique_source_name = df.loc[df['source_center'] == i, 'source_name'].unique()
    if pd.isna(unique_source_name):
        print("Source Center :", i, "-" * 10, "Source Name :", 'Not Found')
    else :
        print("Source Center :", i, "-" * 10, "Source Name :", unique_source_name)
```

```
Source Center : IND342902A1B ---------- Source Name : Not Found
Source Center : IND577116AAA ---------- Source Name : Not Found
Source Center : IND282002AAD ---------- Source Name : Not Found
Source Center : IND465333A1B ---------- Source Name : Not Found
Source Center : IND841301AAC ---------- Source Name : Not Found
Source Center : IND509103AAC ---------- Source Name : Not Found
Source Center : IND126116AAA ---------- Source Name : Not Found
Source Center : IND331022A1B ---------- Source Name : Not Found
Source Center : IND505326AAB ---------- Source Name : Not Found
Source Center : IND852118A1B ---------- Source Name : Not Found
```

```python
for i in missing_source_name:
    unique_destination_name = df.loc[df['destination_center'] == i, 'destination_name'].unique()
    if (pd.isna(unique_source_name)) or (unique_source_name.size == 0):
        print("Destination Center :", i, "-" * 10, "Destination Name :", 'Not Found')
    else :
        print("Destination Center :", i, "-" * 10, "Destination Name :", unique_destination_name)
```

```
Destination Center : IND342902A1B ---------- Destination Name : Not Found
Destination Center : IND577116AAA ---------- Destination Name : Not Found
Destination Center : IND282002AAD ---------- Destination Name : Not Found
Destination Center : IND465333A1B ---------- Destination Name : Not Found
Destination Center : IND841301AAC ---------- Destination Name : Not Found
Destination Center : IND509103AAC ---------- Destination Name : Not Found
Destination Center : IND126116AAA ---------- Destination Name : Not Found
Destination Center : IND331022A1B ---------- Destination Name : Not Found
Destination Center : IND505326AAB ---------- Destination Name : Not Found
Destination Center : IND852118A1B ---------- Destination Name : Not Found
```

```python
missing_destination_name = df.loc[df['destination_name'].isnull(), 'destination_center'].unique()
missing_destination_name
```

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

```python
np.all(df.loc[df['source_name'].isnull(), 'source_center'].isin(missing_destination_name))
```

```
False
```

```python
count = 1
for i in missing_destination_name:
    df.loc[df['destination_center'] == i, 'destination_name'] = df.loc[df['destination_center'] == i, 'destination_name'].replace(np.nan, f'
    count += 1
```

```python
d = {}
for i in missing_source_name:
    d[i] = df.loc[df['destination_center'] == i, 'destination_name'].unique()
for idx, val in d.items():
    if len(val) == 0:
        d[idx] = [f'location_{count}']
        count += 1
d2 = {}
for idx, val in d.items():
    d2[idx] = val[0]
for i, v in d2.items():
    print(i, v)
```

```
IND342902A1B location_1
IND577116AAA location_2
IND282002AAD location_3
IND465333A1B location_4
IND841301AAC location_5
IND509103AAC location_9
IND126116AAA location_8
IND331022A1B location_14
```

```
        IND505326AAB location_6
        IND852118A1B location_7
```

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

```python
df.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
actual_distance_to_destination    0
actual_time                       0
osrm_time                         0
osrm_distance                     0
segment_actual_time               0
segment_osrm_time                 0
segment_osrm_distance             0
dtype: int64
```

```python
df.describe()
```

|       | start_scan_to_end_scan | actual_distance_to_destination | actual_time   | osrm_t     |
|-------|------------------------|--------------------------------|---------------|------------|
| count | 144867.000000          | 144867.000000                  | 144867.000000 | 144867.000 |
| mean  | 961.262986             | 234.073380                     | 416.927521    | 213.868    |
| std   | 1037.012769            | 344.990021                     | 598.103638    | 308.011    |
| min   | 20.000000              | 9.000046                       | 9.000000      | 6.000      |
| 25%   | 161.000000             | 23.355875                      | 51.000000     | 27.000     |
| 50%   | 449.000000             | 66.126572                      | 132.000000    | 64.000     |
| 75%   | 1634.000000            | 286.708878                     | 513.000000    | 257.000    |
| max   | 7898.000000            | 1927.447754                    | 4532.000000   | 1686.000   |

```python
df.describe(include = 'object')
```

|        | route_schedule_uuid                              | trip_uuid              | source_center | source_name                      | de |
|--------|--------------------------------------------------|------------------------|---------------|----------------------------------|----|
| count  | 144867                                           | 144867                 | 144867        | 144867                           |    |
| unique | 1504                                             | 14817                  | 1508          | 1508                             |    |
| top    | thanos::sroute:4029a8a2-6c74-4b7e-a6d8-f9e069f   | trip-153811219535896559 | IND000000ACB  | Gurgaon_Bilaspur_HB (Haryana)    |    |

```python
grouping_1 = ['trip_uuid', 'source_center', 'destination_center']
df1 = df.groupby(by = grouping_1, as_index = False).agg({'data' : 'first',
                                        'route_type' : 'first',
                                        'trip_creation_time' : 'first',
                                        'source_name' : 'first',
                                        'destination_name' : 'last',
                                        'od_start_time' : 'first',
                                        'od_end_time' : 'first',
                                        'start_scan_to_end_scan' : 'first',
                                        'actual_distance_to_destination' : 'last',
                                        'actual_time' : 'last',
                                        'osrm_time' : 'last',
                                        'osrm_distance' : 'last',
                                        'segment_actual_time' : 'sum',
                                        'segment_osrm_time' : 'sum',
                                        'segment_osrm_distance' : 'sum'})
```

```python
df1['od_total_time'] = df1['od_end_time'] - df1['od_start_time']
df1.drop(columns = ['od_end_time', 'od_start_time'], inplace = True)
df1['od_total_time'] = df1['od_total_time'].apply(lambda x : round(x.total_seconds() / 60.0, 2))
df1['od_total_time'].head()
```

```
0    1260.60
1     999.51
2      58.83
3     122.78
4     834.64
Name: od_total_time, dtype: float64
```

```python
df2 = df1.groupby(by = 'trip_uuid', as_index = False).agg({'source_center' : 'first',
                                                            'destination_center' : 'last',
                                                            'data' : 'first',
                                                            'route_type' : 'first',
                                                            'trip_creation_time' : 'first',
                                                            'source_name' : 'first',
                                                            'destination_name' : 'last',
                                                            'od_total_time' : 'sum',
                                                            'start_scan_to_end_scan' : 'sum',
                                                            'actual_distance_to_destination' : 'sum',
                                                            'actual_time' : 'sum',
                                                            'osrm_time' : 'sum',
                                                            'osrm_distance' : 'sum',
                                                            'segment_actual_time' : 'sum',
                                                            'segment_osrm_time' : 'sum',
                                                            'segment_osrm_distance' : 'sum'})
```

```python
def location_name_to_state(x):
    l = x.split('(')
    if len(l) == 1:
        return l[0]
    else:
        return l[1].replace(')', "")


def location_name_to_city(x):
    if 'location' in x:
        return 'unknown_city'
    else:
        l = x.split()[0].split('_')
        if 'CCU' in x:
            return 'Kolkata'
        elif 'MAA' in x.upper():
            return 'Chennai'
        elif ('HBR' in x.upper()) or ('BLR' in x.upper()):
            return 'Bengaluru'
        elif 'FBD' in x.upper():
            return 'Faridabad'
        elif 'BOM' in x.upper():
            return 'Mumbai'
        elif 'DEL' in x.upper():
            return 'Delhi'
        elif 'OK' in x.upper():
            return 'Delhi'
        elif 'GZB' in x.upper():
            return 'Ghaziabad'
        elif 'GGN' in x.upper():
            return 'Gurgaon'
        elif 'AMD' in x.upper():
            return 'Ahmedabad'
        elif 'CJB' in x.upper():
            return 'Coimbatore'
        elif 'HYD' in x.upper():
            return 'Hyderabad'
        return l[0]
```

```python
def location_name_to_place(x):
    if 'location' in x:
        return x
    elif 'HBR' in x:
        return 'HBR Layout PC'
    else:
        l = x.split()[0].split('_', 1)
        if len(l) == 1:
            return 'unknown_place'
        else:
            return l[1]
```

```python
df2['source_state'] = df2['source_name'].apply(location_name_to_state)
df2['source_state'].unique()
```

```
array(['Uttar Pradesh', 'Karnataka', 'Haryana', 'Maharashtra',
       'Tamil Nadu', 'Gujarat', 'Delhi', 'Telangana', 'Rajasthan',
       'Assam', 'Madhya Pradesh', 'West Bengal', 'Andhra Pradesh',
       'Punjab', 'Chandigarh', 'Goa', 'Jharkhand', 'Pondicherry',
       'Orissa', 'Uttarakhand', 'Himachal Pradesh', 'Kerala',
       'Arunachal Pradesh', 'Bihar', 'Chhattisgarh',
       'Dadra and Nagar Haveli', 'Jammu & Kashmir', 'Mizoram', 'Nagaland',
       'location_9', 'location_3', 'location_2', 'location_14',
       'location_7'], dtype=object)
```

```python
df2['source_city'] = df2['source_name'].apply(location_name_to_city)
print('No of source cities :', df2['source_city'].nunique())
df2['source_city'].unique()[:100]
```

```
No of source cities : 690
array(['Kanpur', 'Doddablpur', 'Gurgaon', 'Mumbai', 'Bellary', 'Chennai',
       'Bengaluru', 'Surat', 'Delhi', 'Pune', 'Faridabad', 'Shirala',
       'Hyderabad', 'Thirumalagiri', 'Gulbarga', 'Jaipur', 'Allahabad',
       'Guwahati', 'Narsinghpur', 'Shrirampur', 'Madakasira', 'Sonari',
       'Dindigul', 'Jalandhar', 'Chandigarh', 'Deoli', 'Pandharpur',
       'Kolkata', 'Bhandara', 'Kurnool', 'Bhiwandi', 'Bhatinda',
       'RoopNagar', 'Bantwal', 'Lalru', 'Kadi', 'Shahdol', 'Gangakher',
       'Durgapur', 'Vapi', 'Jamjodhpur', 'Jetpur', 'Mehsana', 'Jabalpur',
       'Junagadh', 'Gundlupet', 'Mysore', 'Goa', 'Bhopal', 'Sonipat',
       'Himmatnagar', 'Jamshedpur', 'Pondicherry', 'Anand', 'Udgir',
       'Nadiad', 'Villupuram', 'Purulia', 'Bhubaneshwar', 'Bamangola',
       'Tiruppattur', 'Kotdwara', 'Medak', 'Bangalore', 'Dhrangadhra',
       'Hospet', 'Ghumarwin', 'Agra', 'Sitapur', 'Canacona', 'Bilimora',
       'SultnBthry', 'Lucknow', 'Vellore', 'Bhuj', 'Dinhata',
       'Margherita', 'Boisar', 'Vizag', 'Tezpur', 'Koduru', 'Tirupati',
       'Pen', 'Ahmedabad', 'Faizabad', 'Gandhinagar', 'Anantapur',
       'Betul', 'Panskura', 'Rasipurm', 'Sankari', 'Jorhat', 'PNQ',
       'Srikakulam', 'Dehradun', 'Jassur', 'Sawantwadi', 'Shajapur',
       'Ludhiana', 'GreaterThane'], dtype=object)
```

```python
df2['source_place'] = df2['source_name'].apply(location_name_to_place)
df2['source_place'].unique()[:100]
```

```
array(['Central_H_6', 'ChikaDPP_D', 'Bilaspur_HB', 'unknown_place', 'Dc',
       'Poonamallee', 'Chrompet_DPC', 'HBR Layout PC', 'Central_D_12',
       'Lajpat_IP', 'North_D_3', 'Balabhgarh_DPC', 'Central_DPP_3',
       'Shamshbd_H', 'Xroad_D', 'Nehrugnj_I', 'Central_I_7',
       'Central_H_1', 'Nangli_IP', 'North', 'KndliDPP_D', 'Central_D_9',
       'DavkharRd_D', 'Bandel_D', 'RTCStand_D', 'Central_DPP_1',
       'KGAirprt_HB', 'North_D_2', 'Central_D_1', 'DC', 'Mthurard_L',
       'Mullanpr_DC', 'Central_DPP_2', 'RajCmplx_D', 'Beliaghata_DPC',
       'RjnaiDPP_D', 'AbbasNgr_I', 'Mankoli_HB', 'DPC', 'Airport_H',
       'Hub', 'Gateway_HB', 'Tathawde_H', 'ChotiHvl_DC', 'Trmltmpl_D',
       'OnkarDPP_D', 'Mehmdpur_H', 'KaranNGR_D', 'Sohagpur_D',
       'Chrompet_L', 'Busstand_D', 'Central_I_1', 'IndEstat_I', 'Court_D',
       'Panchot_IP', 'Adhartal_IP', 'DumDum_DPC', 'Bomsndra_HB',
       'Swamylyt_D', 'Yadvgiri_IP', 'Old', 'Kundli_H', 'Central_I_3',
       'Vasanthm_I', 'Poonamallee_HB', 'VUNagar_DC', 'NlgaonRd_D',
       'Bnnrghta_L', 'Thirumtr_IP', 'GariDPP_D', 'Jogshwri_I',
       'KoilStrt_D', 'CotnGren_M', 'Nzbadrd_D', 'Dwaraka_D', 'Nelmngla_H',
       'NvygRDPP_D', 'Gndhichk_D', 'Central_D_3', 'Chowk_D', 'CharRsta_D',
       'Kollgpra_D', 'Peenya_IP', 'GndhiNgr_IP', 'Sanpada_I',
       'WrdN4DPP_D', 'Sakinaka_RP', 'CivilHPL_D', 'OstwlEmp_D',
       'Gajuwaka', 'Mhbhirab_D', 'MGRoad_D', 'Balajicly_I', 'BljiMrkt_D',
       'Dankuni_HB', 'Trnsport_H', 'Rakhial', 'Memnagar', 'East_I_21',
       'Mithakal_D'], dtype=object)
```

```python
df2['destination_state'] = df2['destination_name'].apply(location_name_to_state)
df2['destination_state'].head(10)
```

```
0     Uttar Pradesh
1         Karnataka
2           Haryana
3       Maharashtra
4         Karnataka
5        Tamil Nadu
6        Tamil Nadu
7         Karnataka
8           Gujarat
9             Delhi
Name: destination_state, dtype: object
```

```python
df2['destination_city'] = df2['destination_name'].apply(location_name_to_city)
df2['destination_city'].head()
```

```
0        Kanpur
1     Doddablpur
2       Gurgaon
3        Mumbai
4        Sandur
Name: destination_city, dtype: object
```

```python
df2['destination_place'] = df2['destination_name'].apply(location_name_to_place)
df2['destination_place'].head()
```

```
0     Central_H_6
1       ChikaDPP_D
2      Bilaspur_HB
3        MiraRd_IP
4       WrdN1DPP_D
Name: destination_place, dtype: object
```

```python
df2['trip_creation_date'] = pd.to_datetime(df2['trip_creation_time'].dt.date)
df2['trip_creation_date'].head()
```

```
0    2018-09-12
1    2018-09-12
2    2018-09-12
3    2018-09-12
4    2018-09-12
Name: trip_creation_date, dtype: datetime64[ns]
```

```python
df2['trip_creation_month'] = df2['trip_creation_time'].dt.month
df2['trip_creation_month'] = df2['trip_creation_month'].astype('int8')
df2['trip_creation_month'].head()
```

```
0    9
1    9
2    9
3    9
4    9
Name: trip_creation_month, dtype: int8
```

```python
df2['trip_creation_year'] = df2['trip_creation_time'].dt.year
df2['trip_creation_year'] = df2['trip_creation_year'].astype('int16')
df2['trip_creation_year'].head()
```

```
0    2018
1    2018
2    2018
3    2018
4    2018
Name: trip_creation_year, dtype: int16
```

```python
df2['trip_creation_week'] = df2['trip_creation_time'].dt.isocalendar().week
df2['trip_creation_week'] = df2['trip_creation_week'].astype('int8')
df2['trip_creation_week'].head()
```

```
0    37
1    37
2    37
3    37
```

```
     4    37
     Name: trip_creation_week, dtype: int8
```

```python
df2['trip_creation_hour'] = df2['trip_creation_time'].dt.hour
df2['trip_creation_hour'] = df2['trip_creation_hour'].astype('int8')
df2['trip_creation_hour'].head()
```

```
     0    0
     1    0
     2    0
     3    0
     4    0
     Name: trip_creation_hour, dtype: int8
```

```python
df2.shape
```

```
     (14817, 28)
```

```python
df2.info()
```

```
     <class 'pandas.core.frame.DataFrame'>
     RangeIndex: 14817 entries, 0 to 14816
     Data columns (total 28 columns):
      #   Column                          Non-Null Count  Dtype
     ---  ------                          --------------  -----
      0   trip_uuid                       14817 non-null  object
      1   source_center                   14817 non-null  object
      2   destination_center              14817 non-null  object
      3   data                            14817 non-null  category
      4   route_type                      14817 non-null  category
      5   trip_creation_time              14817 non-null  datetime64[ns]
      6   source_name                     14817 non-null  object
      7   destination_name                14817 non-null  object
      8   od_total_time                   14817 non-null  float64
      9   start_scan_to_end_scan          14817 non-null  float64
      10  actual_distance_to_destination  14817 non-null  float32
      11  actual_time                     14817 non-null  float32
      12  osrm_time                       14817 non-null  float32
      13  osrm_distance                   14817 non-null  float32
      14  segment_actual_time             14817 non-null  float32
      15  segment_osrm_time               14817 non-null  float32
      16  segment_osrm_distance           14817 non-null  float32
      17  source_state                    14817 non-null  object
      18  source_city                     14817 non-null  object
      19  source_place                    14817 non-null  object
      20  destination_state               14817 non-null  object
      21  destination_city                14817 non-null  object
      22  destination_place               14817 non-null  object
      23  trip_creation_date              14817 non-null  datetime64[ns]
      24  trip_creation_month             14817 non-null  int8
      25  trip_creation_year              14817 non-null  int16
      26  trip_creation_week              14817 non-null  int8
      27  trip_creation_hour              14817 non-null  int8
     dtypes: category(2), datetime64[ns](2), float32(7), float64(2), int16(1), int8(3), object(11)
     memory usage: 2.2+ MB
```

```python
df2.describe().T
```

| | count | mean | std | min | 25% |
|---|---|---|---|---|---|
| od_total_time | 14817.0 | 531.697630 | 658.868223 | 23.460000 | 149.930000 |
| start_scan_to_end_scan | 14817.0 | 530.810016 | 658.705957 | 23.000000 | 149.000000 |
| actual_distance_to_destination | 14817.0 | 164.477829 | 305.388153 | 9.002461 | 22.837238 |
| actual_time | 14817.0 | 357.143768 | 561.396118 | 9.000000 | 67.000000 |
| osrm_time | 14817.0 | 161.384018 | 271.360992 | 6.000000 | 29.000000 |
| osrm_distance | 14817.0 | 204.344711 | 370.395569 | 9.072900 | 30.819201 |
| segment_actual_time | 14817.0 | 353.892273 | 556.247925 | 9.000000 | 66.000000 |
| segment_osrm_time | 14817.0 | 180.949783 | 314.542053 | 6.000000 | 31.000000 |
| segment_osrm_distance | 14817.0 | 223.201157 | 416.628387 | 9.072900 | 32.654499 |
| trip_creation_month | 14817.0 | 9.120672 | 0.325757 | 9.000000 | 9.000000 |
| trip_creation_year | 14817.0 | 2018.000000 | 0.000000 | 2018.000000 | 2018.000000 |
| trip_creation_week | 14817.0 | 38.295944 | 0.967872 | 37.000000 | 38.000000 |
| trip_creation_hour | 14817.0 | 12.449821 | 7.986553 | 0.000000 | 4.000000 |

```
df2.describe(include = object).T
```

| | count | unique | top | freq |
|---|---|---|---|---|
| trip_uuid | 14817 | 14817 | trip-153671041653548748 | 1 |
| source_center | 14817 | 938 | IND000000ACB | 1063 |
| destination_center | 14817 | 1042 | IND000000ACB | 821 |
| source_name | 14817 | 938 | Gurgaon_Bilaspur_HB (Haryana) | 1063 |
| destination_name | 14817 | 1042 | Gurgaon_Bilaspur_HB (Haryana) | 821 |
| source_state | 14817 | 34 | Maharashtra | 2714 |
| source_city | 14817 | 690 | Mumbai | 1442 |
| source_place | 14817 | 761 | Bilaspur_HB | 1063 |
| destination_state | 14817 | 39 | Maharashtra | 2561 |
| destination_city | 14817 | 806 | Mumbai | 1548 |
| destination_place | 14817 | 850 | Bilaspur_HB | 821 |

```
df2['trip_creation_hour'].unique()
```

```
array([ 0,  1,  2,  3,  4,  5,  6,  7,  8,  9, 10, 11, 12, 13, 14, 15, 16,
       17, 18, 19, 20, 21, 22, 23], dtype=int8)
```

```
df_hour = df2.groupby(by = 'trip_creation_hour')['trip_uuid'].count().to_frame().reset_index()
df_hour.head()
```

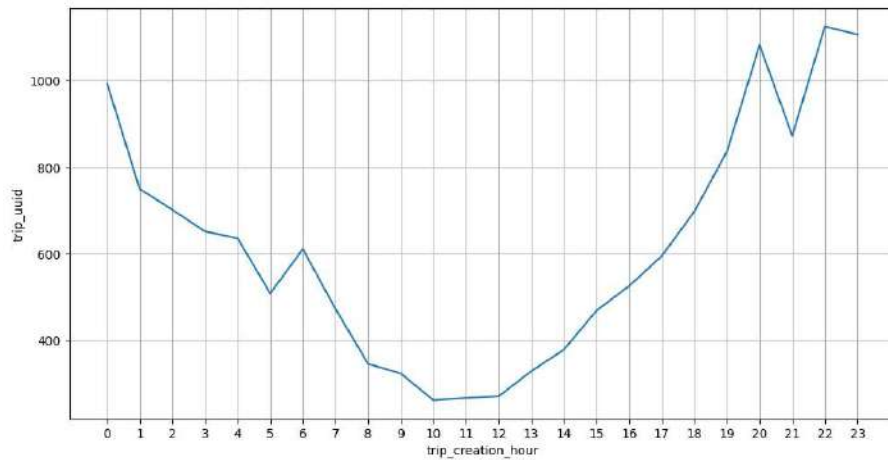| | trip_creation_hour | trip_uuid |
|---|---|---|
| 0 | 0 | 994 |
| 1 | 1 | 750 |
| 2 | 2 | 702 |
| 3 | 3 | 652 |
| 4 | 4 | 636 |

Next steps:  [ Generate code with df_hour ]   [ ◯ View recommended plots ]

```
plt.figure(figsize = (12, 6))
sns.lineplot(data = df_hour,
        x = df_hour['trip_creation_hour'],
        y = df_hour['trip_uuid'],
        markers = '*')
plt.xticks(np.arange(0,24))
```

```python
plt.grid('both')
plt.plot()
```

[]



```python
df2['trip_creation_week'].unique()
```

```
array([37, 38, 39, 40], dtype=int8)
```

```python
df_week = df2.groupby(by = 'trip_creation_week')['trip_uuid'].count().to_frame().reset_index()
df_week.head()
```
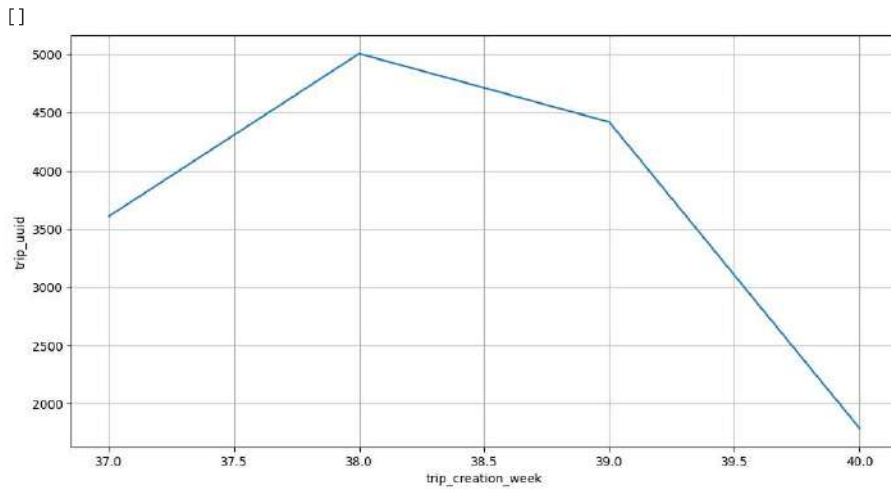
|   | trip_creation_week | trip_uuid |
|---|---|---|
| 0 | 37 | 3608 |
| 1 | 38 | 5004 |
| 2 | 39 | 4417 |
| 3 | 40 | 1788 |

Next steps:    Generate code with `df_week`        ◉ View recommended plots

```python
plt.figure(figsize = (12, 6))
sns.lineplot(data = df_week,
             x = df_week['trip_creation_week'],
             y = df_week['trip_uuid'],
             markers = 'o')
plt.grid('both')
plt.plot()
```

[]



```python
df_month = df2.groupby(by = 'trip_creation_month')['trip_uuid'].count().to_frame().reset_index()
df_month['perc'] = np.round(df_month['trip_uuid'] * 100/ df_month['trip_uuid'].sum(), 2)
df_month.head()
```
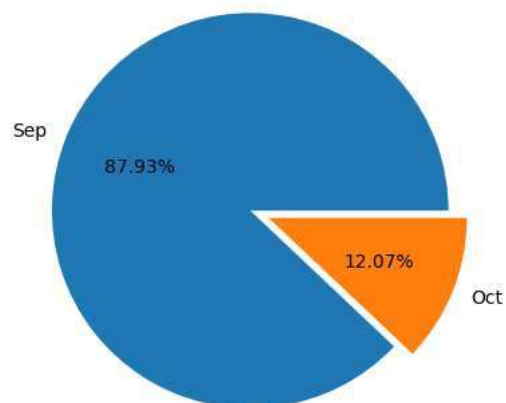
| | trip_creation_month | trip_uuid | perc |
|---|---|---|---|
| 0 | 9 | 13029 | 87.93 |
| 1 | 10 | 1788 | 12.07 |

Next steps:    Generate code with `df_month`        ◉ View recommended plots

```python
plt.pie(x = df_month['trip_uuid'],
        labels = ['Sep', 'Oct'],
        explode = [0, 0.1],
       autopct = '%.2f%%')
plt.plot()
```

[]



```python
df_data = df2.groupby(by = 'data')['trip_uuid'].count().to_frame().reset_index()
df_data['perc'] = np.round(df_data['trip_uuid'] * 100/ df_data['trip_uuid'].sum(), 2)
df_data.head()
```
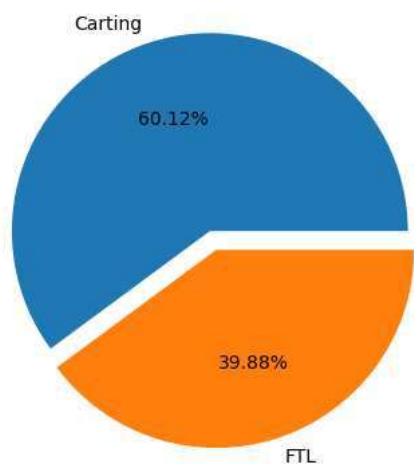
|   | data | trip_uuid | perc |
|---|------|-----------|------|
| 0 | test | 4163 | 28.1 |
| 1 | training | 10654 | 71.9 |

Next steps:  [ Generate code with `df_data` ]    [⊙ View recommended plots ]

```
plt.pie(x = df_data['trip_uuid'],
        labels = df_data['data'],
        explode = [0, 0.1],
        autopct = '%.2f%%')
plt.plot()
```

    []



```
df_route = df2.groupby(by = 'route_type')['trip_uuid'].count().to_frame().reset_index()
df_route['perc'] = np.round(df_route['trip_uuid'] * 100/ df_route['trip_uuid'].sum(), 2)
df_route.head()
```

|   | route_type | trip_uuid | perc |
|---|------------|-----------|------|
| 0 | Carting | 8908 | 60.12 |
| 1 | FTL | 5909 | 39.88 |

Next steps:  [ Generate code with `df_route` ]    [⊙ View recommended plots ]

```
plt.pie(x = df_route['trip_uuid'],
        labels = ['Carting', 'FTL'],
        explode = [0, 0.1],
        autopct = '%.2f%%')
plt.plot()
```

[]



```python
df_source_state = df2.groupby(by = 'source_state')['trip_uuid'].count().to_frame().reset_index()
df_source_state['perc'] = np.round(df_source_state['trip_uuid'] * 100/ df_source_state['trip_uuid'].sum(), 2)
df_source_state = df_source_state.sort_values(by = 'trip_uuid', ascending = False)
df_source_state.head()
```
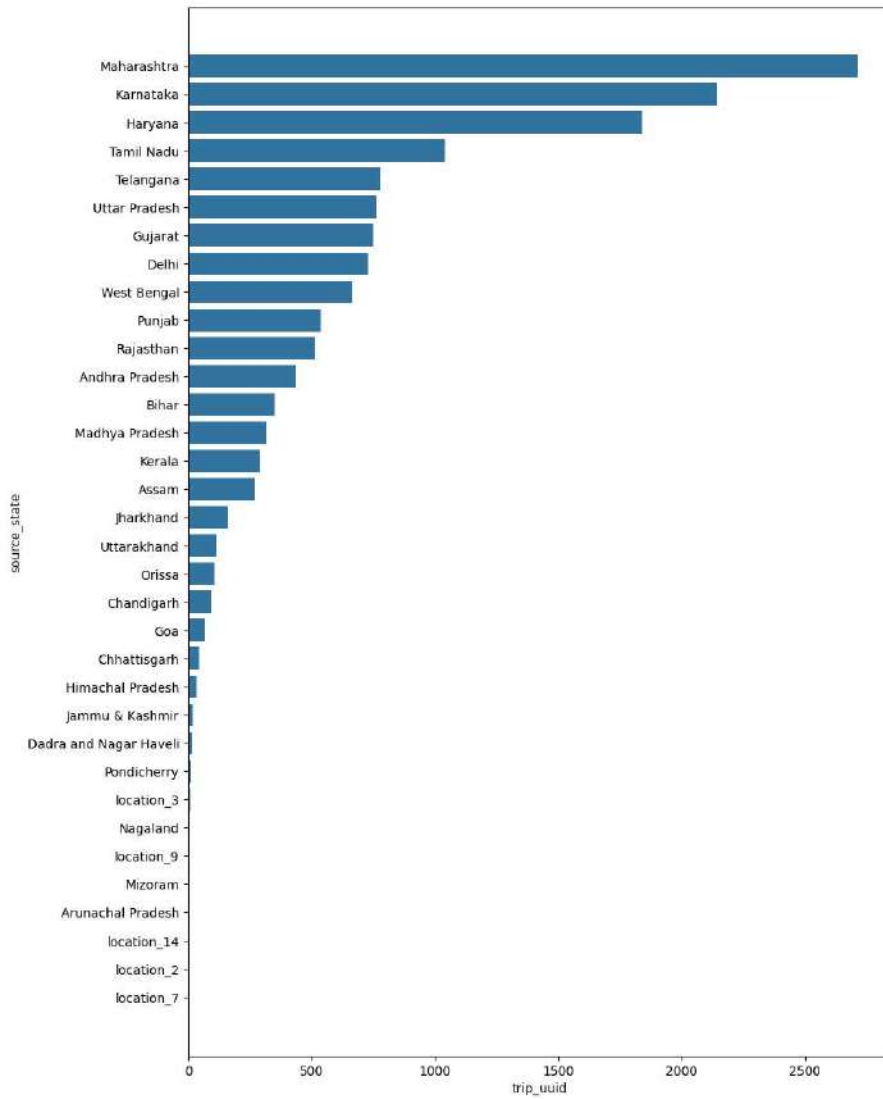
|    | source_state | trip_uuid | perc  |
|----|--------------|-----------|-------|
| 17 | Maharashtra  | 2714      | 18.32 |
| 14 | Karnataka    | 2143      | 14.46 |
| 10 | Haryana      | 1838      | 12.40 |
| 24 | Tamil Nadu   | 1039      | 7.01  |
| 25 | Telangana    | 781       | 5.27  |

Next steps:    Generate code with `df_source_state`    ⬤ View recommended plots

```python
plt.figure(figsize = (10, 15))
sns.barplot(data = df_source_state,
            x = df_source_state['trip_uuid'],
            y = df_source_state['source_state'])
plt.plot()
```

[]



```python
df_source_city = df2.groupby(by = 'source_city')['trip_uuid'].count().to_frame().reset_index()
df_source_city['perc'] = np.round(df_source_city['trip_uuid'] * 100/ df_source_city['trip_uuid'].sum(), 2)
df_source_city = df_source_city.sort_values(by = 'trip_uuid', ascending = False)[:30]
df_source_city
```
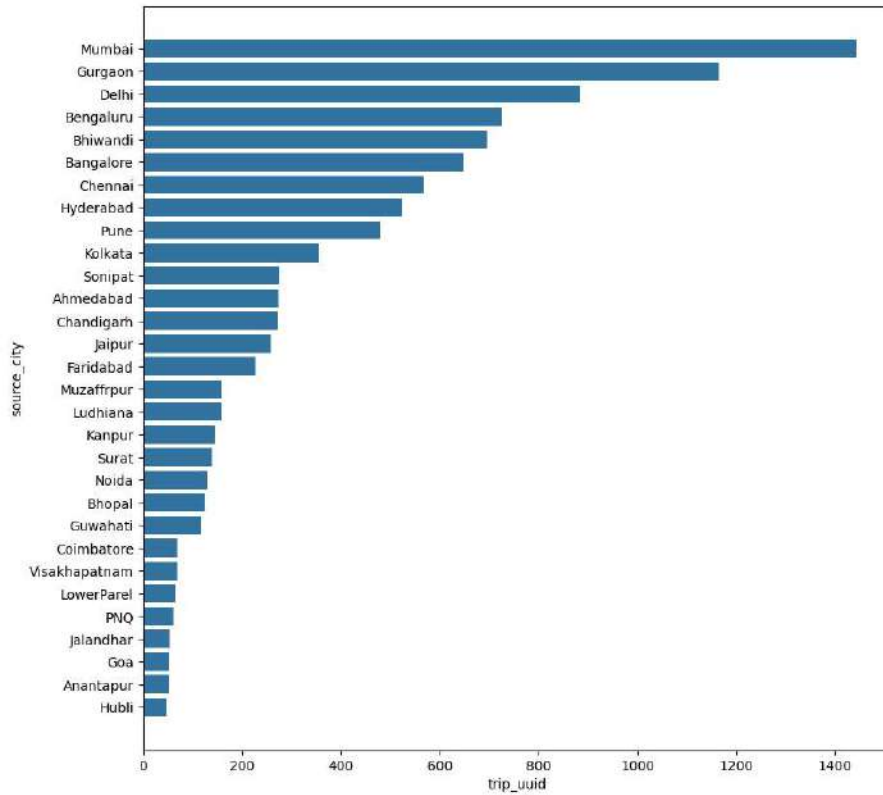
|     | source_city   | trip_uuid | perc |
|-----|---------------|-----------|------|
| 439 | Mumbai        | 1442      | 9.73 |
| 237 | Gurgaon       | 1165      | 7.86 |
| 169 | Delhi         | 883       | 5.96 |
| 79  | Bengaluru     | 726       | 4.90 |
| 100 | Bhiwandi      | 697       | 4.70 |
| 58  | Bangalore     | 648       | 4.37 |
| 136 | Chennai       | 568       | 3.83 |
| 264 | Hyderabad     | 524       | 3.54 |
| 516 | Pune          | 480       | 3.24 |
| 357 | Kolkata       | 356       | 2.40 |
| 610 | Sonipat       | 276       | 1.86 |
| 2   | Ahmedabad     | 274       | 1.85 |
| 133 | Chandigarh    | 273       | 1.84 |
| 270 | Jaipur        | 259       | 1.75 |
| 201 | Faridabad     | 227       | 1.53 |
| 447 | Muzaffrpur    | 159       | 1.07 |
| 382 | Ludhiana      | 158       | 1.07 |
| 320 | Kanpur        | 145       | 0.98 |
| 621 | Surat         | 140       | 0.94 |
| 473 | Noida         | 129       | 0.87 |
| 102 | Bhopal        | 125       | 0.84 |
| 240 | Guwahati      | 118       | 0.80 |
| 154 | Coimbatore    | 69        | 0.47 |
| 679 | Visakhapatnam | 69        | 0.47 |
| 380 | LowerParel    | 65        | 0.44 |
| 477 | PNQ           | 62        | 0.42 |
| 273 | Jalandhar     | 54        | 0.36 |
| 220 | Goa           | 52        | 0.35 |
| 25  | Anantapur     | 51        | 0.34 |
| 261 | Hubli         | 47        | 0.32 |

Next steps:     Generate code with `df_source_city`     View recommended plots

```
plt.figure(figsize = (10, 10))
sns.barplot(data = df_source_city,
          x = df_source_city['trip_uuid'],
          y = df_source_city['source_city'])
plt.plot()
```
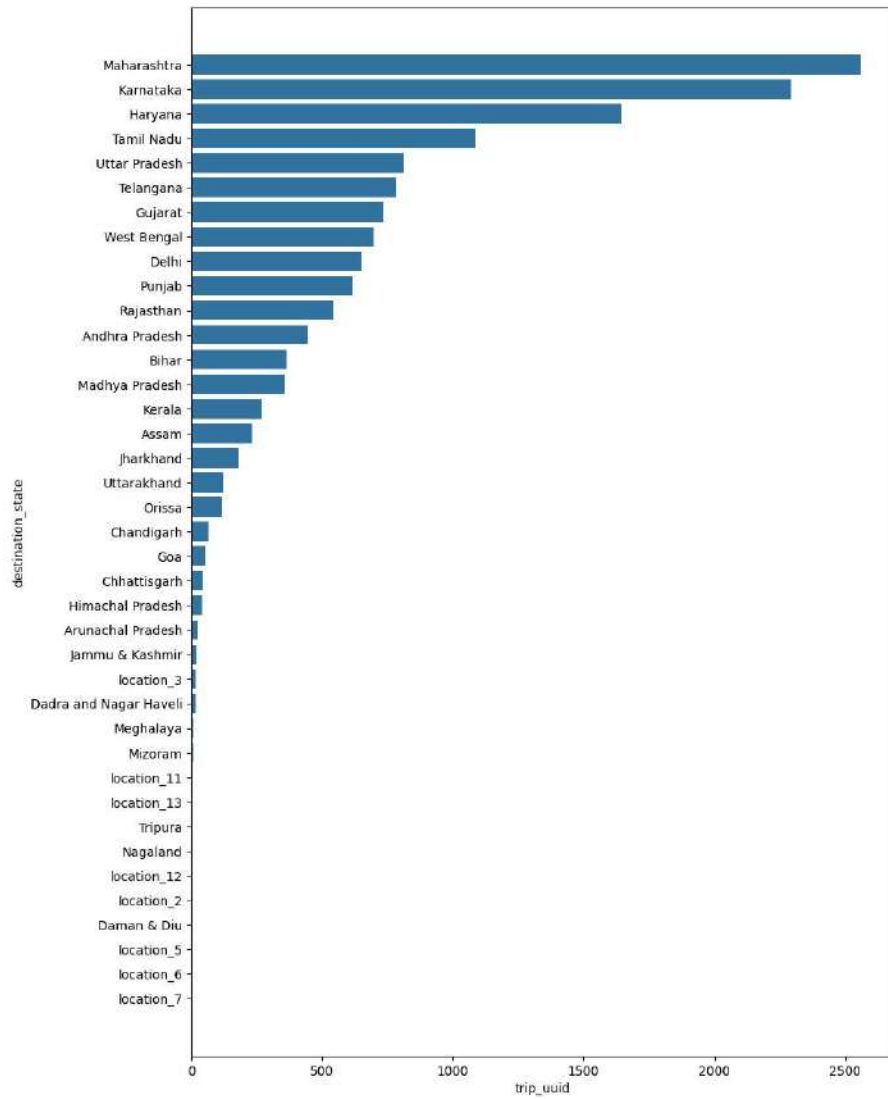
[]



```python
df_destination_state = df2.groupby(by = 'destination_state')['trip_uuid'].count().to_frame().reset_index()
df_destination_state['perc'] = np.round(df_destination_state['trip_uuid'] * 100/ df_destination_state['trip_uuid'].sum(), 2)
df_destination_state = df_destination_state.sort_values(by = 'trip_uuid', ascending = False)
df_destination_state.head()
```

|    | destination_state | trip_uuid | perc  |
|----|-------------------|-----------|-------|
| 18 | Maharashtra       | 2561      | 17.28 |
| 15 | Karnataka         | 2294      | 15.48 |
| 11 | Haryana           | 1643      | 11.09 |
| 25 | Tamil Nadu        | 1084      | 7.32  |
| 28 | Uttar Pradesh     | 811       | 5.47  |

Next steps:     Generate code with `df_destination_state`     View recommended plots

```python
plt.figure(figsize = (10, 15))
sns.barplot(data = df_destination_state,
            x = df_destination_state['trip_uuid'],
            y = df_destination_state['destination_state'])
plt.plot()
```

[]



```python
df_destination_city = df2.groupby(by = 'destination_city')['trip_uuid'].count().to_frame().reset_index()
df_destination_city['perc'] = np.round(df_destination_city['trip_uuid'] * 100/ df_destination_city['trip_uuid'].sum(), 2)
df_destination_city = df_destination_city.sort_values(by = 'trip_uuid', ascending = False)[:30]
df_destination_city
```
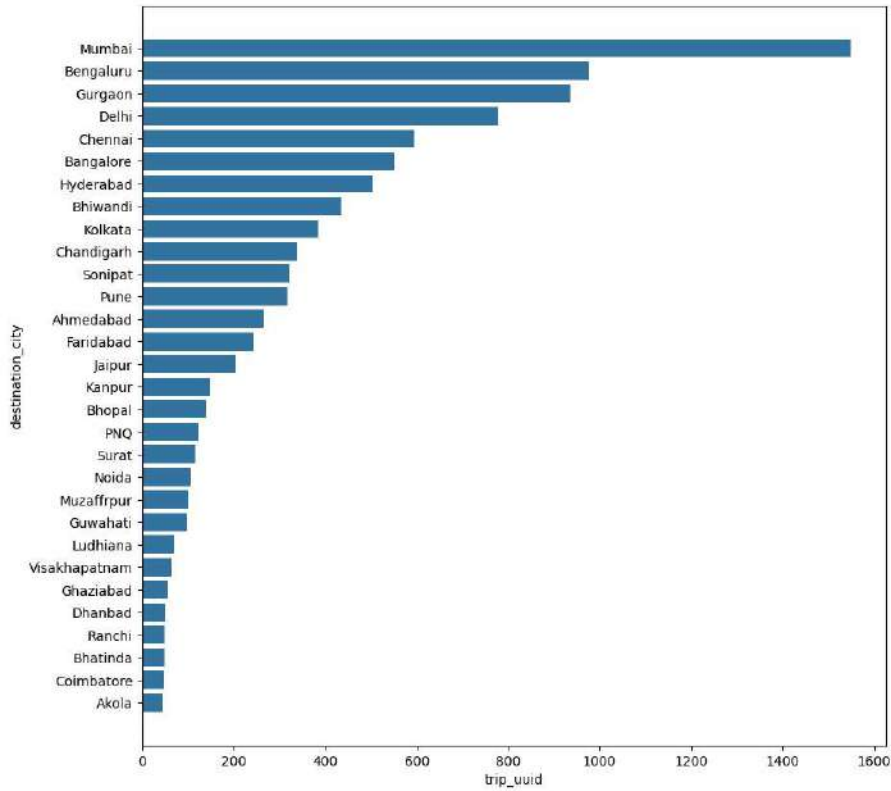
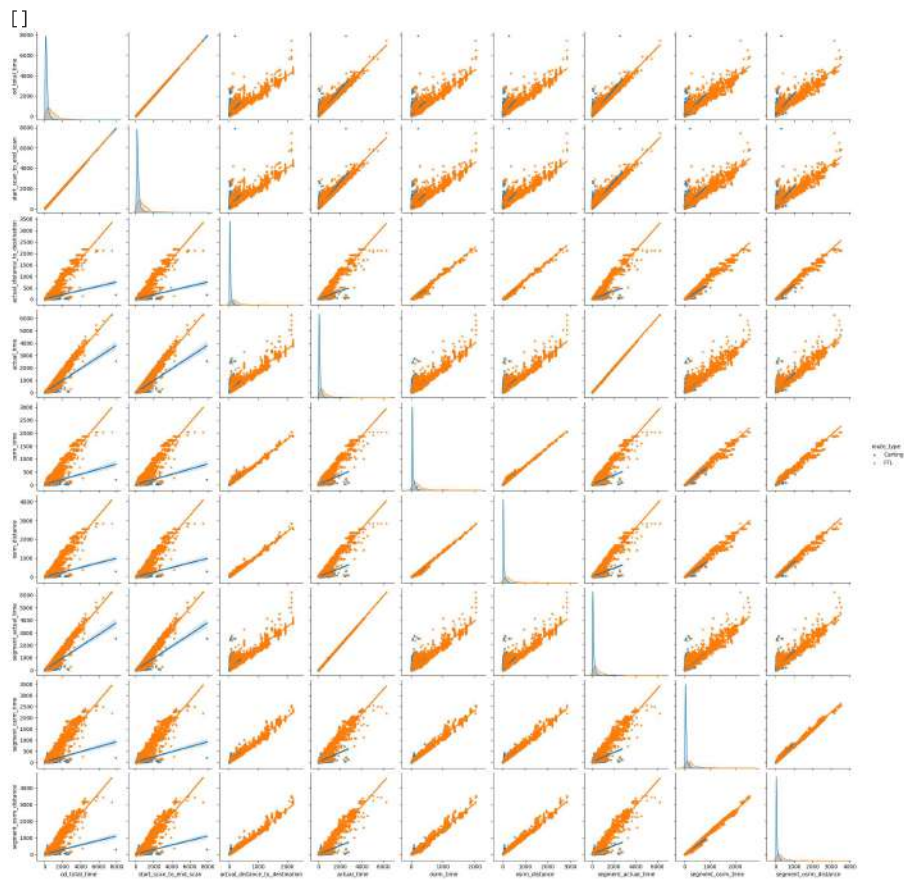|     | destination_city | trip_uuid | perc  |
| --- | --- | --- | --- |
| 515 | Mumbai | 1548 | 10.45 |
| 96  | Bengaluru | 975 | 6.58 |
| 282 | Gurgaon | 936 | 6.32 |
| 200 | Delhi | 778 | 5.25 |
| 163 | Chennai | 595 | 4.02 |
| 72  | Bangalore | 551 | 3.72 |
| 308 | Hyderabad | 503 | 3.39 |
| 115 | Bhiwandi | 434 | 2.93 |
| 418 | Kolkata | 384 | 2.59 |
| 158 | Chandigarh | 339 | 2.29 |
| 724 | Sonipat | 322 | 2.17 |
| 612 | Pune | 317 | 2.14 |
| 4   | Ahmedabad | 265 | 1.79 |
| 242 | Faridabad | 244 | 1.65 |
| 318 | Jaipur | 205 | 1.38 |
| 371 | Kanpur | 148 | 1.00 |
| 117 | Bhopal | 139 | 0.94 |
| 559 | PNQ | 122 | 0.82 |
| 739 | Surat | 117 | 0.79 |
| 552 | Noida | 106 | 0.72 |
| 521 | Muzaffrpur | 102 | 0.69 |
| 284 | Guwahati | 98 | 0.66 |
| 448 | Ludhiana | 70 | 0.47 |
| 797 | Visakhapatnam | 64 | 0.43 |
| 259 | Ghaziabad | 56 | 0.38 |
| 208 | Dhanbad | 50 | 0.34 |
| 639 | Ranchi | 49 | 0.33 |
| 110 | Bhatinda | 48 | 0.32 |
| 183 | Coimbatore | 47 | 0.32 |
| 9   | Akola | 45 | 0.30 |

Next steps:     Generate code with `df_destination_city`          View recommended plots

```python
plt.figure(figsize = (10, 10))
sns.barplot(data = df_destination_city,
            x = df_destination_city['trip_uuid'],
            y = df_destination_city['destination_city'])
plt.plot()
```

[]



```
numerical_columns = ['od_total_time', 'start_scan_to_end_scan', 'actual_distance_to_destination',
                     'actual_time', 'osrm_time', 'osrm_distance', 'segment_actual_time',
                     'segment_osrm_time', 'segment_osrm_distance']
sns.pairplot(data = df2,
            vars = numerical_columns,
            kind = 'reg',
            hue = 'route_type',
            markers = '.')
plt.plot()
```

[]



```
df_corr = df2[numerical_columns].corr()
df_corr
```

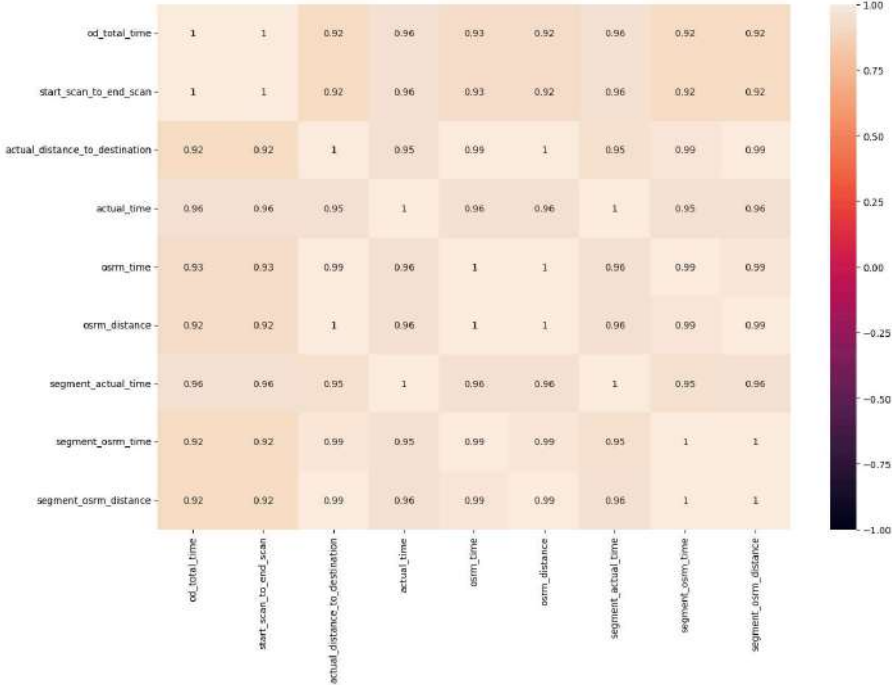|  | od_total_time | start_scan_to_end_scan | actual_distance_to_ |
|---|---|---|---|
| od_total_time | 1.000000 | 0.999999 | |
| start_scan_to_end_scan | 0.999999 | 1.000000 | |
| actual_distance_to_destination | 0.918222 | 0.918308 | |
| actual_time | 0.961094 | 0.961147 | |
| osrm_time | 0.926516 | 0.926571 | |
| osrm_distance | 0.924219 | 0.924299 | |
| segment_actual_time | 0.961119 | 0.961171 | |
| segment_osrm_time | 0.918490 | 0.918561 | |
| segment_osrm_distance | 0.919199 | 0.919291 | |

Next steps:    Generate code with `df_corr`        ◉ View recommended plots

```
plt.figure(figsize = (15, 10))
sns.heatmap(data = df_corr, vmin = -1, vmax = 1, annot = True)
plt.plot()
```

[]



```
df2[['od_total_time', 'start_scan_to_end_scan']].describe()
```

|       | od_total_time | start_scan_to_end_scan |
|-------|---------------|------------------------|
| count | 14817.000000  | 14817.000000           |
| mean  | 531.697630    | 530.810016             |
| std   | 658.868223    | 658.705957             |
| min   | 23.460000     | 23.000000              |
| 25%   | 149.930000    | 149.000000             |
| 50%   | 280.770000    | 280.000000             |
| 75%   | 638.200000    | 637.000000             |
| max   | 7898.550000   | 7898.000000            |

```
plt.figure(figsize = (12, 6))
sns.histplot(df2['od_total_time'], element = 'step', color = 'green')
```

# Business Insights

- The data is given from the period '2018-09-12 00:00:16' to '2018-10-08 03:00:24'.

- There are about 14817 unique trip IDs, 1508 unique source centers, 1481 unique destination_centers, 690 unique source cities, 806 unique destination cities.

- Most of the data is for testing than for training.

- Most common route type is Carting.

- The names of 14 unique location ids are missing in the data.

- The number of trips start increasing after the noon, becomes maximum at 10 P.M and then start decreasing.

- Maximum trips are created in the 38th week.

- Most orders come mid-month. That means customers usually make more orders in the mid of the month.

- Most orders are sourced from the states like Maharashtra, Karnataka, Haryana, Tamil Nadu, Telangana

- Maximum number of trips originated from Mumbai city followed by Gurgaon Delhi, Bengaluru and Bhiwandi. That means that the seller base is strong in these cities.

- Maximum number of trips ended in Maharashtra state followed by Karnataka, Haryana, Tamil Nadu and Uttar Pradesh. That means that the number of orders placed in these states is significantly high.

- Maximum number of trips ended in Mumbai city followed by Bengaluru, Gurgaon, Delhi and Chennai. That means that the number of orders placed in these cities is significantly high.

# Recommendations

- The OSRM trip planning system needs to be improved. Discrepancies need to be catered to for transporters, if the routing engine is configured for optimum results.

- osrm_time and actual_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 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.

- Most of the orders are coming from/reaching to states like Maharashtra, Karnataka, Haryana and Tamil Nadu. The existing corridors can be further enhanced to improve the penetration in these areas.

- Customer profiling of the customers belonging to the states Maharashtra, Karnataka, Haryana, Tamil Nadu and Uttar Pradesh has to be done to get to know why major orders are coming from these atates and to improve customers' buying and delivery experience.

- From state point of view, we might have very heavy traffic in certain states and bad terrain conditions in certain states. This will be a good indicator to plan and cater to demand during peak festival seasons.