

Delhivery - Feature Engineering : Case Study

by Pavan Kumar

About Delhivery

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.

The Data team builds intelligence and capabilities using this data that helps them to widen the gap between the quality, efficiency, and profitability of their business versus their competitors.

How can you help here?

The company wants to understand and process the data coming out of data engineering pipelines:

- Clean, sanitize and manipulate data to get useful features out of raw fields
- Make sense out of the raw data and help the data science team to build forecasting models on it

Column Profiling:

data - tells whether the data is testing or training data trip_creation_time – Timestamp of trip creation route_schedule_uuid – Unique Id for a particular route schedule route_type – Transportation type FTL – Full Truck Load: FTL shipments get to the destination sooner, as the truck is making no other pickups or drop-offs along the way Carting: Handling system consisting of small vehicles (carts) trip_uuid - Unique ID given to a particular trip (A trip may include different source and destination centers) source_center - Source ID of trip origin source_name - Source Name of trip origin destination_center – Destination ID destination_name – Destination Name od_start_time – Trip start time od_end_time – Trip end time start_scan_to_end_scan – Time taken to deliver from source to destination is_cutoff – Unknown field cutoff_factor – Unknown field cutoff_timestamp – Unknown field actual_distance_to_destination – Distance in Kms between source and destination warehouse actual_time – Actual time taken to complete the delivery (Cumulative) osrm_time – An open-source routing engine time calculator which computes the shortest

path between points in a given map (Includes usual traffic, distance through major and minor roads) and gives the time (Cumulative)
osrm_distance – An open-source routing engine which computes the shortest path between points in a given map (Includes usual traffic, distance through major and minor roads) (Cumulative) factor – Unknown field segment_actual_time – This is a segment time. Time taken by the subset of the package delivery segment_osrm_time – This is the OSRM segment time. Time taken by the subset of the package delivery segment_osrm_distance – This is the OSRM distance. Distance covered by subset of the package delivery segment_factor – Unknown field
Concept Used:

Feature Creation

Relationship between Features Column Normalization /Column Standardization Handling categorical values Missing values - Outlier treatment / Types of outliers

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
```

```
In [2]: df = pd.read_csv("delhivery_data.csv")
```

```
In [3]: df.head()
```

Out[3]:

	data	trip_creation_time	route_schedule_uuid	route_type	trip_uuid	source_center	source_name	destination_center	desti
0	training	2018-09-20 02:35:36.476840	thanos::sroute:eb7bfc78- b351-4c0e-a951- fa3d5c3...	Carting	trip- 153741093647649320	IND388121AAA	Anand_VUNagar_DC (Gujarat)	IND388620AAB	Khambhat_
1	training	2018-09-20 02:35:36.476840	thanos::sroute:eb7bfc78- b351-4c0e-a951- fa3d5c3...	Carting	trip- 153741093647649320	IND388121AAA	Anand_VUNagar_DC (Gujarat)	IND388620AAB	Khambhat_
2	training	2018-09-20 02:35:36.476840	thanos::sroute:eb7bfc78- b351-4c0e-a951- fa3d5c3...	Carting	trip- 153741093647649320	IND388121AAA	Anand_VUNagar_DC (Gujarat)	IND388620AAB	Khambhat_
3	training	2018-09-20 02:35:36.476840	thanos::sroute:eb7bfc78- b351-4c0e-a951- fa3d5c3...	Carting	trip- 153741093647649320	IND388121AAA	Anand_VUNagar_DC (Gujarat)	IND388620AAB	Khambhat_
4	training	2018-09-20 02:35:36.476840	thanos::sroute:eb7bfc78- b351-4c0e-a951- fa3d5c3...	Carting	trip- 153741093647649320	IND388121AAA	Anand_VUNagar_DC (Gujarat)	IND388620AAB	Khambhat_

5 rows × 24 columns



In [4]: `df.info()`

```

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

```

```
In [5]: df.isna().sum()
```

```
Out[5]: data          0
trip_creation_time    0
route_schedule_uuid   0
route_type            0
trip_uuid             0
source_center         0
source_name           293
destination_center     0
destination_name       261
od_start_time         0
od_end_time           0
start_scan_to_end_scan 0
is_cutoff             0
cutoff_factor         0
cutoff_timestamp      0
actual_distance_to_destination 0
actual_time           0
osrm_time             0
osrm_distance         0
factor               0
segment_actual_time   0
segment_osrm_time     0
segment_osrm_distance 0
segment_factor        0
dtype: int64
```

```
In [6]: df.isnull().sum()
```

```
Out[6]: data
trip_creation_time      0
route_schedule_uuid     0
route_type              0
trip_uuid               0
source_center           0
source_name             293
destination_center      0
destination_name        261
od_start_time           0
od_end_time             0
start_scan_to_end_scan  0
is_cutoff               0
cutoff_factor           0
cutoff_timestamp        0
actual_distance_to_destination 0
actual_time             0
osrm_time               0
osrm_distance           0
factor                  0
segment_actual_time     0
segment_osrm_time       0
segment_osrm_distance   0
segment_factor          0
dtype: int64
```

- Removing Null values

```
In [7]: df = df.dropna(how='any')
df = df.reset_index(drop=True)
```

```
In [8]: df.isna().sum()
```

```
Out[8]: data          0
trip_creation_time    0
route_schedule_uuid   0
route_type            0
trip_uuid             0
source_center         0
source_name           0
destination_center    0
destination_name      0
od_start_time         0
od_end_time           0
start_scan_to_end_scan 0
is_cutoff             0
cutoff_factor         0
cutoff_timestamp      0
actual_distance_to_destination 0
actual_time           0
osrm_time             0
osrm_distance         0
factor               0
segment_actual_time   0
segment_osrm_time     0
segment_osrm_distance 0
segment_factor        0
dtype: int64
```

```
In [9]: df.isnull().sum()
```

```
Out[9]: data          0
trip_creation_time    0
route_schedule_uuid   0
route_type            0
trip_uuid             0
source_center         0
source_name           0
destination_center    0
destination_name      0
od_start_time         0
od_end_time           0
start_scan_to_end_scan 0
is_cutoff             0
cutoff_factor         0
cutoff_timestamp      0
actual_distance_to_destination 0
actual_time           0
osrm_time             0
osrm_distance         0
factor               0
segment_actual_time   0
segment_osrm_time     0
segment_osrm_distance 0
segment_factor        0
dtype: int64
```

```
In [10]: df.describe()
```


Out[10]:

	start_scan_to_end_scan	cutoff_factor	actual_distance_to_destination	actual_time	osrm_time	osrm_distance	factor	segment_actual_
count	144316.000000	144316.000000	144316.000000	144316.000000	144316.000000	144316.000000	144316.000000	144316.00
mean	963.697698	233.561345	234.708498	417.996237	214.437055	285.549785	2.120178	36.17
std	1038.082976	345.245823	345.480571	598.940065	308.448543	421.717826	1.717065	53.52
min	20.000000	9.000000	9.000045	9.000000	6.000000	9.008200	0.144000	-244.00
25%	161.000000	22.000000	23.352027	51.000000	27.000000	29.896250	1.604545	20.00
50%	451.000000	66.000000	66.135322	132.000000	64.000000	78.624400	1.857143	28.00
75%	1645.000000	286.000000	286.919294	516.000000	259.000000	346.305400	2.212280	40.00
max	7898.000000	1927.000000	1927.447705	4532.000000	1686.000000	2326.199100	77.387097	3051.00

- Convert date time

```
In [11]: df['od_start_time'] = pd.to_datetime(df['od_start_time'])  
df['od_end_time'] = pd.to_datetime(df['od_end_time'])
```

```
In [12]: df.head()
```

Out[12]:

	data	trip_creation_time	route_schedule_uuid	route_type	trip_uuid	source_center	source_name	destination_center	desti
0	training	2018-09-20 02:35:36.476840	thanos::sroute:eb7bfc78- b351-4c0e-a951- fa3d5c3...	Carting	153741093647649320	IND388121AAA	Anand_VUNagar_DC (Gujarat)	IND388620AAB	Khambhat_
1	training	2018-09-20 02:35:36.476840	thanos::sroute:eb7bfc78- b351-4c0e-a951- fa3d5c3...	Carting	153741093647649320	IND388121AAA	Anand_VUNagar_DC (Gujarat)	IND388620AAB	Khambhat_
2	training	2018-09-20 02:35:36.476840	thanos::sroute:eb7bfc78- b351-4c0e-a951- fa3d5c3...	Carting	153741093647649320	IND388121AAA	Anand_VUNagar_DC (Gujarat)	IND388620AAB	Khambhat_
3	training	2018-09-20 02:35:36.476840	thanos::sroute:eb7bfc78- b351-4c0e-a951- fa3d5c3...	Carting	153741093647649320	IND388121AAA	Anand_VUNagar_DC (Gujarat)	IND388620AAB	Khambhat_
4	training	2018-09-20 02:35:36.476840	thanos::sroute:eb7bfc78- b351-4c0e-a951- fa3d5c3...	Carting	153741093647649320	IND388121AAA	Anand_VUNagar_DC (Gujarat)	IND388620AAB	Khambhat_

5 rows × 24 columns

- Grouping by sub-journey in trip

```
In [13]: df['segment_key'] = df['trip_uuid'] + df['source_center'] + df['destination_center']

segment_cols = ['segment_actual_time', 'segment_osrm_distance', 'segment_osrm_time']

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

df[[col + '_sum' for col in segment_cols]]
```

Out[13]:

	segment_actual_time_sum	segment_osrm_distance_sum	segment_osrm_time_sum
0	14.0	11.9653	11.0
1	24.0	21.7243	20.0
2	40.0	32.5395	27.0
3	61.0	45.5619	39.0
4	67.0	49.4772	44.0
...
144311	92.0	65.3487	94.0
144312	118.0	82.7212	115.0
144313	138.0	103.4265	149.0
144314	155.0	122.3150	176.0
144315	423.0	131.1238	185.0

144316 rows × 3 columns

- aggregating at sub-journey

In [14]:

```
create_segment_dict = {  
    'data' : 'first',  
    'trip_creation_time' : 'first',  
    'route_schedule_uuid' : 'first',  
    'route_type' : 'first',  
    'trip_uuid' : 'first',  
    'source_center' : 'first',  
    'source_name' : 'first',  
  
    'destination_center' : 'last',  
    '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' : 'last',  
'segment_osrm_distance_sum' : 'last',  
'segment_osrm_time_sum' : 'last',  
  
}
```

- Grouping mini-trips, sorting by time

```
In [15]: segment = df.groupby('segment_key').agg(create_segment_dict).reset_index()  
segment = segment.sort_values(by=['segment_key', 'od_end_time'], ascending=True).reset_index()
```

```
In [16]: segment
```

Out[16]:

	index	segment_key	data	trip_creation_time	route_schedule_uuid	route_type	trip_uuid	sol
0	0	trip-153671041653548748IND209304AAAIND000000ACB	training	2018-09-12 00:00:16.535741	thanos::sroute:d7c989ba-a29b-4a0b-b2f4-288cdc6...	FTL	trip-153671041653548748	IND
1	1	trip-153671041653548748IND462022AAAIND209304AAA	training	2018-09-12 00:00:16.535741	thanos::sroute:d7c989ba-a29b-4a0b-b2f4-288cdc6...	FTL	trip-153671041653548748	IND
2	2	trip-153671042288605164IND561203AABIND562101AAA	training	2018-09-12 00:00:22.886430	thanos::sroute:3a1b0ab2-bb0b-4c53-8c59-eb2a2c0...	Carting	trip-153671042288605164	IND
3	3	trip-153671042288605164IND572101AAAIND561203AAB	training	2018-09-12 00:00:22.886430	thanos::sroute:3a1b0ab2-bb0b-4c53-8c59-eb2a2c0...	Carting	trip-153671042288605164	IND
4	4	trip-153671043369099517IND000000ACBIND160002AAC	training	2018-09-12 00:00:33.691250	thanos::sroute:de5e208e-7641-45e6-8100-4d9fb1e...	FTL	trip-153671043369099517	IND
...
26217	26217	trip-153861115439069069IND628204AAAIND627657AAA	test	2018-10-03 23:59:14.390954	thanos::sroute:c5f2ba2c-8486-4940-8af6-d1d2a6a...	Carting	trip-153861115439069069	IND
26218	26218	trip-153861115439069069IND628613AAAIND627005AAA	test	2018-10-03 23:59:14.390954	thanos::sroute:c5f2ba2c-8486-4940-8af6-d1d2a6a...	Carting	trip-153861115439069069	IND
26219	26219	trip-153861115439069069IND628801AAAIND628204AAA	test	2018-10-03 23:59:14.390954	thanos::sroute:c5f2ba2c-8486-4940-8af6-d1d2a6a...	Carting	trip-153861115439069069	IND
26220	26220	trip-153861118270144424IND583119AAAIND583101AAA	test	2018-10-03 23:59:42.701692	thanos::sroute:412fea14-6d1f-4222-8a5f-a517042...	FTL	trip-153861118270144424	IND
26221	26221	trip-153861118270144424IND583201AAAIND583119AAA	test	2018-10-03 23:59:42.701692	thanos::sroute:412fea14-6d1f-4222-8a5f-a517042...	FTL	trip-153861118270144424	IND

26222 rows × 21 columns

```
In [17]: segment[segment['trip_uuid'] == 'trip-153741093647649320' ]
```

Out[17]:

	index		segment_key	data	trip_creation_time	route_schedule_uuid	route_type	trip_uuid	source	
10370	10370	153741093647649320	trip-IND388121AAA	IND388620AAB	training	2018-09-20 02:35:36.476840	thanos::sroute:eb7bfc78-b351-4c0e-a951-fa3d5c3...	Carting	153741093647649320	IND388121AAA
10371	10371	153741093647649320	trip-IND388620AAB	IND388320AAA	training	2018-09-20 02:35:36.476840	thanos::sroute:eb7bfc78-b351-4c0e-a951-fa3d5c3...	Carting	153741093647649320	IND388620AAB

2 rows × 21 columns



```
In [18]: segment.info()
```



```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 26222 entries, 0 to 26221
Data columns (total 21 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   index                                26222 non-null  int64
1   segment_key                          26222 non-null  object
2   data                                26222 non-null  object
3   trip_creation_time                   26222 non-null  object
4   route_schedule_uuid                  26222 non-null  object
5   route_type                           26222 non-null  object
6   trip_uuid                            26222 non-null  object
7   source_center                        26222 non-null  object
8   source_name                          26222 non-null  object
9   destination_center                   26222 non-null  object
10  destination_name                      26222 non-null  object
11  od_start_time                         26222 non-null  datetime64[ns]
12  od_end_time                           26222 non-null  datetime64[ns]
13  start_scan_to_end_scan                26222 non-null  float64
14  actual_distance_to_destination         26222 non-null  float64
15  actual_time                           26222 non-null  float64
16  osrm_time                             26222 non-null  float64
17  osrm_distance                         26222 non-null  float64
18  segment_actual_time_sum                26222 non-null  float64
19  segment_osrm_distance_sum              26222 non-null  float64
20  segment_osrm_time_sum                  26222 non-null  float64
dtypes: datetime64[ns](2), float64(8), int64(1), object(10)
memory usage: 4.2+ MB

```

Consider time taken between `od_start_time` and `od_end_time` as a feature.

```

In [19]: segment['od_time_diff_hour'] = (segment['od_end_time'] - segment['od_start_time']).dt.total_seconds() / (60)
segment['od_time_diff_hour']

```

```
Out[19]: 0      1260.604421
          1      999.505379
          2      58.832388
          3     122.779486
          4     834.638929
          ...
          26217    62.115193
          26218    91.087797
          26219    44.174403
          26220   287.474007
          26221    66.933565
          Name: od_time_diff_hour, Length: 26222, dtype: float64
```

```
In [20]: segment
```


Out[20]:

	index	segment_key	data	trip_creation_time	route_schedule_uuid	route_type	trip_uuid	sol
0	0	153671041653548748IND209304AAAIND000000ACB	trip- training	2018-09-12 00:00:16.535741	thanos::sroute:d7c989ba-a29b-4a0b-b2f4-288cdc6...	FTL	trip- 153671041653548748	IND
1	1	153671041653548748IND462022AAAIND209304AAA	trip- training	2018-09-12 00:00:16.535741	thanos::sroute:d7c989ba-a29b-4a0b-b2f4-288cdc6...	FTL	trip- 153671041653548748	IND
2	2	153671042288605164IND561203AABIND562101AAA	trip- training	2018-09-12 00:00:22.886430	thanos::sroute:3a1b0ab2-bb0b-4c53-8c59-eb2a2c0...	Carting	trip- 153671042288605164	IND
3	3	153671042288605164IND572101AAAIND561203AAB	trip- training	2018-09-12 00:00:22.886430	thanos::sroute:3a1b0ab2-bb0b-4c53-8c59-eb2a2c0...	Carting	trip- 153671042288605164	IND
4	4	153671043369099517IND000000ACBIND160002AAC	trip- training	2018-09-12 00:00:33.691250	thanos::sroute:de5e208e-7641-45e6-8100-4d9fb1e...	FTL	trip- 153671043369099517	IND
...
26217	26217	153861115439069069IND628204AAAIND627657AAA	trip- test	2018-10-03 23:59:14.390954	thanos::sroute:c5f2ba2c-8486-4940-8af6-d1d2a6a...	Carting	trip- 153861115439069069	IND
26218	26218	153861115439069069IND628613AAAIND627005AAA	trip- test	2018-10-03 23:59:14.390954	thanos::sroute:c5f2ba2c-8486-4940-8af6-d1d2a6a...	Carting	trip- 153861115439069069	IND
26219	26219	153861115439069069IND628801AAAIND628204AAA	trip- test	2018-10-03 23:59:14.390954	thanos::sroute:c5f2ba2c-8486-4940-8af6-d1d2a6a...	Carting	trip- 153861115439069069	IND
26220	26220	153861118270144424IND583119AAAIND583101AAA	trip- test	2018-10-03 23:59:42.701692	thanos::sroute:412fea14-6d1f-4222-8a5f-a517042...	FTL	trip- 153861118270144424	IND
26221	26221	153861118270144424IND583201AAAIND583119AAA	trip- test	2018-10-03 23:59:42.701692	thanos::sroute:412fea14-6d1f-4222-8a5f-a517042...	FTL	trip- 153861118270144424	IND

26222 rows × 22 columns

```
In [21]: create_trip_dict = {  
    'data' : 'first',  
    'trip_creation_time' : 'first',  
    'route_schedule_uuid' : 'first',  
    'route_type' : 'first',  
    'trip_uuid' : 'first',  
  
    'source_center' : 'first',  
    'source_name' : 'first',  
  
    'destination_center' : 'last',  
    'destination_name' : 'last',  
  
    'start_scan_to_end_scan' : 'sum',  
    'od_time_diff_hour' : 'sum',  
  
    'actual_distance_to_destination' : 'sum',  
    'actual_time' : 'sum',  
    'osrm_time' : 'sum',  
    'osrm_distance' : 'sum',  
  
    'segment_actual_time_sum' : 'sum',  
    'segment_osrm_distance_sum' : 'sum',  
    'segment_osrm_time_sum' : 'sum',  
  
}
```

```
In [22]: trip = segment.groupby('trip_uuid').agg(create_trip_dict).reset_index(drop = True)
```

```
In [23]: trip
```

Out[23]:

	data	trip_creation_time	route_schedule_uuid	route_type	trip_uuid	source_center	source_name	destination_center
0	training	2018-09-12 00:00:16.535741	thanos::sroute:d7c989ba- a29b-4a0b-b2f4- 288cdc6...	FTL	trip- 153671041653548748	IND209304AAA	Kanpur_Central_H_6 (Uttar Pradesh)	IND209304AAA
1	training	2018-09-12 00:00:22.886430	thanos::sroute:3a1b0ab2- bb0b-4c53-8c59- eb2a2c0...	Carting	trip- 153671042288605164	IND561203AAB	Doddablpur_ChikaDPP_D (Karnataka)	IND561203AAB
2	training	2018-09-12 00:00:33.691250	thanos::sroute:de5e208e- 7641-45e6-8100- 4d9fb1e...	FTL	trip- 153671043369099517	IND000000ACB	Gurgaon_Bilaspur_HB (Haryana)	IND000000ACB
3	training	2018-09-12 00:01:00.113710	thanos::sroute:f0176492- a679-4597-8332- bbd1c7f...	Carting	trip- 153671046011330457	IND400072AAB	Mumbai Hub (Maharashtra)	IND401104AAA
4	training	2018-09-12 00:02:09.740725	thanos::sroute:d9f07b12- 65e0-4f3b-bec8- df06134...	FTL	trip- 153671052974046625	IND583101AAA	Bellary_Dc (Karnataka)	IND583119AAA
...
14782	test	2018-10-03 23:55:56.258533	thanos::sroute:8a120994- f577-4491-9e4b- b7e4a14...	Carting	trip- 153861095625827784	IND160002AAC	Chandigarh_Mehmdpur_H (Punjab)	IND160002AAC
14783	test	2018-10-03 23:57:23.863155	thanos::sroute:b30e1ec3- 3bfa-4bd2-a7fb- 3b75769...	Carting	trip- 153861104386292051	IND121004AAB	FBD_Balabhgarh_DPC (Haryana)	IND121004AAA
14784	test	2018-10-03 23:57:44.429324	thanos::sroute:5609c268- e436-4e0a-8180- 3db4a74...	Carting	trip- 153861106442901555	IND208006AAA	Kanpur_GovndNgr_DC (Uttar Pradesh)	IND208006AAA
14785	test	2018-10-03 23:59:14.390954	thanos::sroute:c5f2ba2c- 8486-4940-8af6- d1d2a6a...	Carting	trip- 153861115439069069	IND627005AAA	Tirunelveli_VdkkuSrt_I (Tamil Nadu)	IND628204AAA
14786	test	2018-10-03 23:59:42.701692	thanos::sroute:412fea14- 6d1f-4222-8a5f- a517042...	FTL	trip- 153861118270144424	IND583119AAA	Sandur_WrdN1DPP_D (Karnataka)	IND583119AAA

14787 rows × 18 columns

```
In [24]: trip[['actual_time', 'segment_actual_time_sum']]
```

```
Out[24]:
```

	actual_time	segment_actual_time_sum
0	1562.0	1548.0
1	143.0	141.0
2	3347.0	3308.0
3	59.0	59.0
4	341.0	340.0
...
14782	83.0	82.0
14783	21.0	21.0
14784	282.0	281.0
14785	264.0	258.0
14786	275.0	274.0

14787 rows × 2 columns

```
In [25]: trip
```

Out[25]:

	data	trip_creation_time	route_schedule_uuid	route_type	trip_uuid	source_center	source_name	destination_center
0	training	2018-09-12 00:00:16.535741	thanos::sroute:d7c989ba- a29b-4a0b-b2f4- 288cdc6...	FTL	trip- 153671041653548748	IND209304AAA	Kanpur_Central_H_6 (Uttar Pradesh)	IND209304AAA
1	training	2018-09-12 00:00:22.886430	thanos::sroute:3a1b0ab2- bb0b-4c53-8c59- eb2a2c0...	Carting	trip- 153671042288605164	IND561203AAB	Doddablpur_ChikaDPP_D (Karnataka)	IND561203AAB
2	training	2018-09-12 00:00:33.691250	thanos::sroute:de5e208e- 7641-45e6-8100- 4d9fb1e...	FTL	trip- 153671043369099517	IND000000ACB	Gurgaon_Bilaspur_HB (Haryana)	IND000000ACB
3	training	2018-09-12 00:01:00.113710	thanos::sroute:f0176492- a679-4597-8332- bbd1c7f...	Carting	trip- 153671046011330457	IND400072AAB	Mumbai Hub (Maharashtra)	IND401104AAA
4	training	2018-09-12 00:02:09.740725	thanos::sroute:d9f07b12- 65e0-4f3b-bec8- df06134...	FTL	trip- 153671052974046625	IND583101AAA	Bellary_Dc (Karnataka)	IND583119AAA
...
14782	test	2018-10-03 23:55:56.258533	thanos::sroute:8a120994- f577-4491-9e4b- b7e4a14...	Carting	trip- 153861095625827784	IND160002AAC	Chandigarh_Mehmdpur_H (Punjab)	IND160002AAC
14783	test	2018-10-03 23:57:23.863155	thanos::sroute:b30e1ec3- 3bfa-4bd2-a7fb- 3b75769...	Carting	trip- 153861104386292051	IND121004AAB	FBD_Balabhgarh_DPC (Haryana)	IND121004AAA
14784	test	2018-10-03 23:57:44.429324	thanos::sroute:5609c268- e436-4e0a-8180- 3db4a74...	Carting	trip- 153861106442901555	IND208006AAA	Kanpur_GovndNgr_DC (Uttar Pradesh)	IND208006AAA
14785	test	2018-10-03 23:59:14.390954	thanos::sroute:c5f2ba2c- 8486-4940-8af6- d1d2a6a...	Carting	trip- 153861115439069069	IND627005AAA	Tirunelveli_VdkkuSrt_I (Tamil Nadu)	IND628204AAA
14786	test	2018-10-03 23:59:42.701692	thanos::sroute:412fea14- 6d1f-4222-8a5f- a517042...	FTL	trip- 153861118270144424	IND583119AAA	Sandur_WrdN1DPP_D (Karnataka)	IND583119AAA

14787 rows × 18 columns



```
In [26]: trip[trip['trip_uuid'] == 'trip-153741093647649320' ]
```

```
Out[26]:
```

	data	trip_creation_time	route_schedule_uuid	route_type	trip_uuid	source_center	source_name	destination_center	destina
5917	training	2018-09-20 02:35:36.476840	thanos::sroute:eb7bfc78- b351-4c0e-a951- fa3d5c3...	Carting	trip- 153741093647649320	IND388121AAA	Anand_VUNagar_DC (Gujarat)	IND388320AAA	Anand

```
In [27]: trip[['actual_distance_to_destination', 'osrm_distance']]
```

```
Out[27]:
```

	actual_distance_to_destination	osrm_distance
0	824.732854	991.3523
1	73.186911	85.1110
2	1927.404273	2354.0665
3	17.175274	19.6800
4	127.448500	146.7918
...
14782	57.762332	73.4630
14783	15.513784	16.0882
14784	38.684839	58.9037
14785	134.723836	171.1103
14786	66.081533	80.5787

14787 rows × 2 columns

Hypothesis Testing

```
In [28]: trip['destination_name'] = trip['destination_name'].str.lower() # Lowering all columns
trip['source_name'] = trip['source_name']
```

```
In [29]: def place2state(x):
    # transform "gurgaon_bilaspur_hb (haryana)" into "haryana"
    state = x.split('(')[1]

    return state[:-1] #removing ')' from ending

def place2city(x):
    #we will remove state
    city = x.split(' (')[0]

    city = city.split('_')[0]

    # Now daling with edge cases

    if city == 'pnq vadgaon sheri dpc': return 'vadgaonsheri'

    # ['PNQ Pashan DPC', 'Bhopal MP Nagar', 'HBR Layout PC',
    #  'PNQ Rahatani DPC', 'Pune Balaji Nagar', 'Mumbai Antop Hill']

    if city in ['pnq pashan dpc', 'pnq rahatani dpc', 'pune balaji nagar']:
        return 'pune'

    if city == 'hbr layout pc' :
        return 'bengaluru'
    if city == 'bhopal mp nagar':
        return 'bhopal'
    if city == 'mumbai antop hill':
        return 'mumbai'

    return city

def place2city_place(x):

    # we will remove state
    x = x.split('(')[0]

    len_ = len(x.split('_'))

    if len_ >= 3:
        return x.split('_')[1]
```

```

    # small cities have same city and place name
    if len_ == 2:
        return x.split('_')[0]

    # now we need to deal with edge cases or improper name convention

    # if len(x.split('_')) == 2:

    return x.split(' ')[0]

def place2code(x):
    # we will remove state
    x = x.split('(')[0]

    if len(x.split('_')) >= 3:
        return x.split('_')[-1]

    return 'none'

```

```

In [30]: trip['destination_state'] = trip['destination_name'].apply(lambda x: place2state(x))
trip['destination_city'] = trip['destination_name'].apply(lambda x: place2city(x))
trip['destination_place'] = trip['destination_name'].apply(lambda x: place2city_place(x))
trip['destination_code'] = trip['destination_name'].apply(lambda x: place2code(x))

```

```

In [31]: trip[['destination_state', 'destination_city', 'destination_place', 'destination_code']]

```


Out[31]:

	destination_state	destination_city	destination_place	destination_code
0	uttar pradesh	kanpur	central	6
1	karnataka	doddablpur	chikadpp	d
2	haryana	gurgaon	bilaspur	hb
3	maharashtra	mumbai	mirard	ip
4	karnataka	sandur	wrdn1dpp	d
...
14782	punjab	chandigarh	mehmdpur	h
14783	haryana	faridabad	blbgarh	dc
14784	uttar pradesh	kanpur	govndngr	dc
14785	tamil nadu	tirchchnr	shnmgprm	d
14786	karnataka	sandur	wrdn1dpp	d

14787 rows × 4 columns

```
In [32]: trip['trip_creation_time'] = pd.to_datetime(trip['trip_creation_time'])
```

```
trip['trip_year'] = trip['trip_creation_time'].dt.year
trip['trip_month'] = trip['trip_creation_time'].dt.month
trip['trip_hour'] = trip['trip_creation_time'].dt.hour
trip['trip_day'] = trip['trip_creation_time'].dt.day
trip['trip_week'] = trip['trip_creation_time'].dt.isocalendar().week
trip['trip_dayofweek'] = trip['trip_creation_time'].dt.dayofweek
```

```
In [33]: trip[['trip_year','trip_month','trip_hour','trip_day','trip_week','trip_dayofweek']]
```

Out[33]:

	trip_year	trip_month	trip_hour	trip_day	trip_week	trip_dayofweek
0	2018	9	0	12	37	2
1	2018	9	0	12	37	2
2	2018	9	0	12	37	2
3	2018	9	0	12	37	2
4	2018	9	0	12	37	2
...
14782	2018	10	23	3	40	2
14783	2018	10	23	3	40	2
14784	2018	10	23	3	40	2
14785	2018	10	23	3	40	2
14786	2018	10	23	3	40	2

14787 rows × 6 columns

In [34]:

```
trip.head()
```

Out[34]:		data	trip_creation_time	route_schedule_uuid	route_type	trip_uuid	source_center	source_name	destination_center	
0	training		2018-09-12 00:00:16.535741	thanos::sroute:d7c989ba- a29b-4a0b-b2f4- 288cdc6...	FTL	153671041653548748	IND209304AAA	Kanpur_Central_H_6 (Uttar Pradesh)	IND209304AAA	k
1	training		2018-09-12 00:00:22.886430	thanos::sroute:3a1b0ab2- bb0b-4c53-8c59- eb2a2c0...	Carting	153671042288605164	IND561203AAB	Doddablpur_ChikaDPP_D (Karnataka)	IND561203AAB	dodda
2	training		2018-09-12 00:00:33.691250	thanos::sroute:de5e208e- 7641-45e6-8100- 4d9fb1e...	FTL	153671043369099517	IND000000ACB	Gurgaon_Bilaspur_HB (Haryana)	IND000000ACB	gu
3	training		2018-09-12 00:01:00.113710	thanos::sroute:f0176492- a679-4597-8332- bbd1c7f...	Carting	153671046011330457	IND400072AAB	Mumbai Hub (Maharashtra)	IND401104AAA	i
4	training		2018-09-12 00:02:09.740725	thanos::sroute:d9f07b12- 65e0-4f3b-bec8- df06134...	FTL	153671052974046625	IND583101AAA	Bellary_Dc (Karnataka)	IND583119AAA	sa

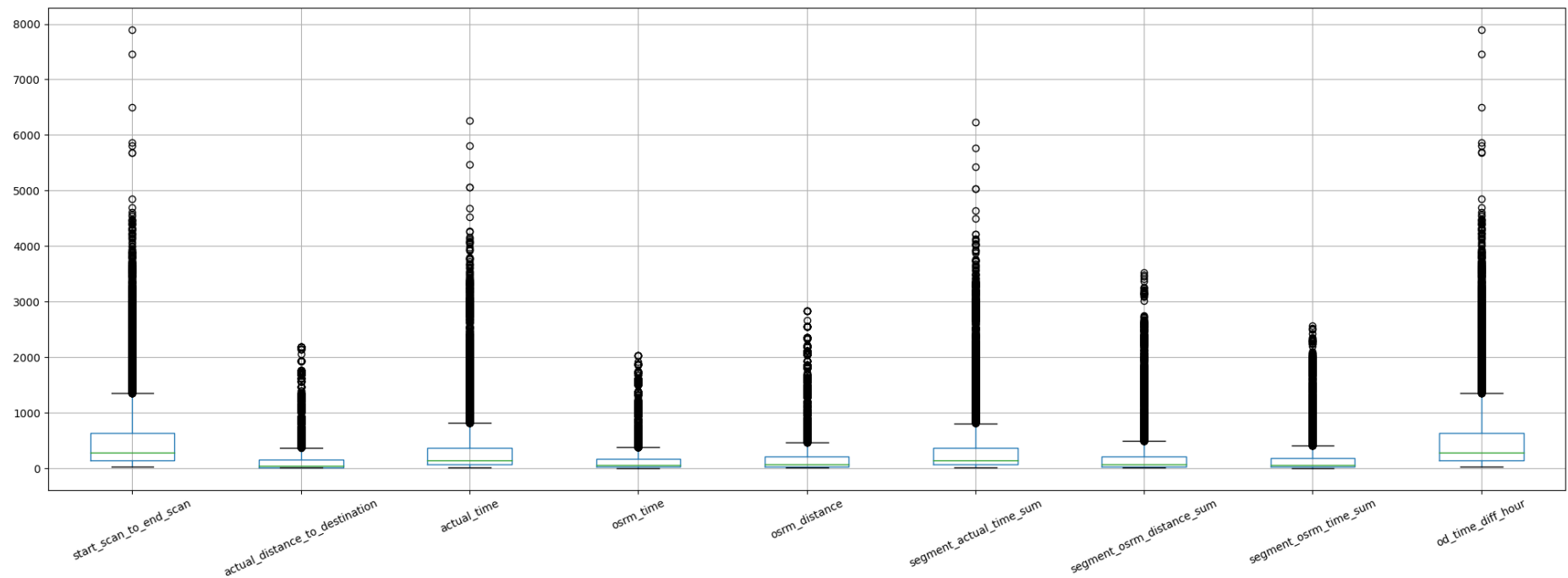
5 rows × 28 columns

```
In [35]: num_cols = ['start_scan_to_end_scan', 'actual_distance_to_destination', 'actual_time', 'osrm_time',
                    'osrm_distance', 'segment_actual_time_sum', 'segment_osrm_distance_sum',
                    'segment_osrm_time_sum', 'od_time_diff_hour']
```

Analysing outliers with box plot

```
In [36]: trip[num_cols].boxplot(rot=25, figsize=(25,8))
```

Out[36]: <Axes: >



Use IQR method to handle the outliers

```
In [37]: Q1 = trip[num_cols].quantile(0.25)
Q3 = trip[num_cols].quantile(0.75)

IQR = Q3 - Q1
```

```
In [38]: trip = trip[~((trip[num_cols] < (Q1 - 1.5 * IQR)) | (trip[num_cols] > (Q3 + 1.5 * IQR))).any(axis=1)]
trip = trip.reset_index(drop=True)
```

```
In [39]: trip
```

Out[39]:

	data	trip_creation_time	route_schedule_uuid	route_type	trip_uuid	source_center	source_name	destination_center
0	training	2018-09-12 00:00:22.886430	thanos::sroute:3a1b0ab2- bb0b-4c53-8c59- eb2a2c0...	Carting	trip- 153671042288605164	IND561203AAB	Doddablpur_ChikaDPP_D (Karnataka)	IND561203AAB
1	training	2018-09-12 00:01:00.113710	thanos::sroute:f0176492- a679-4597-8332- bbd1c7f...	Carting	trip- 153671046011330457	IND400072AAB	Mumbai Hub (Maharashtra)	IND401104AAA
2	training	2018-09-12 00:02:09.740725	thanos::sroute:d9f07b12- 65e0-4f3b-bec8- df06134...	FTL	trip- 153671052974046625	IND583101AAA	Bellary_Dc (Karnataka)	IND583119AAA
3	training	2018-09-12 00:02:34.161600	thanos::sroute:9bf03170- d0a2-4a3f-aa4d- 9aaab3d...	Carting	trip- 153671055416136166	IND600056AAA	Chennai_Poonamallee (Tamil Nadu)	IND600056AAA
4	training	2018-09-12 00:04:22.011653	thanos::sroute:a97698cc- 846e-41a7-916b- 88b1741...	Carting	trip- 153671066201138152	IND600044AAD	Chennai_Chrompet_DPC (Tamil Nadu)	IND600048AAA
...
12718	test	2018-10-03 23:55:56.258533	thanos::sroute:8a120994- f577-4491-9e4b- b7e4a14...	Carting	trip- 153861095625827784	IND160002AAC	Chandigarh_Mehmdpur_H (Punjab)	IND160002AAC
12719	test	2018-10-03 23:57:23.863155	thanos::sroute:b30e1ec3- 3bfa-4bd2-a7fb- 3b75769...	Carting	trip- 153861104386292051	IND121004AAB	FBD_Balabhgarh_DPC (Haryana)	IND121004AAA
12720	test	2018-10-03 23:57:44.429324	thanos::sroute:5609c268- e436-4e0a-8180- 3db4a74...	Carting	trip- 153861106442901555	IND208006AAA	Kanpur_GovndNgr_DC (Uttar Pradesh)	IND208006AAA
12721	test	2018-10-03 23:59:14.390954	thanos::sroute:c5f2ba2c- 8486-4940-8af6- d1d2a6a...	Carting	trip- 153861115439069069	IND627005AAA	Tirunelveli_VdkkuSrt_I (Tamil Nadu)	IND628204AAA
12722	test	2018-10-03 23:59:42.701692	thanos::sroute:412fea14- 6d1f-4222-8a5f- a517042...	FTL	trip- 153861118270144424	IND583119AAA	Sandur_WrdN1DPP_D (Karnataka)	IND583119AAA

12723 rows × 28 columns

```
In [40]: trip[num_cols].boxplot(rot=25, figsize=(25,8))
```

```
Out[40]: <Axes: >
```

Encoding route type

```
In [41]: trip['route_type'].value_counts()
```

```
Out[41]: Carting      8812  
FTL          3911  
Name: route_type, dtype: int64
```

```
In [42]: trip['route_type'] = trip['route_type'].map({'FTL':0, 'Carting':1})
```

Normalize/Standardize the numerical features using MinMaxScaler

```
In [43]: from sklearn.preprocessing import StandardScaler
```

```
In [44]: scaler = StandardScaler()  
scaler.fit(trip[num_cols])
```

```
Out[44]: StandardScaler()
```

```
In [45]: trip[num_cols] = scaler.transform(trip[num_cols])
```

```
In [46]: trip[num_cols]
```

Out[46]:

	start_scan_to_end_scan	actual_distance_to_destination	actual_time	osrm_time	osrm_distance	segment_actual_time_sum	segment_osrm_distance_su
0	-0.548546	0.012060	-0.217856	-0.144341	-0.073948	-0.221500	-0.14531
1	-0.861602	-0.765152	-0.749015	-0.877085	-0.804506	-0.743482	-0.82361
2	1.552838	0.764988	1.034163	0.533102	0.614738	1.045260	0.51489
3	-0.513328	-0.662169	-0.736369	-0.766482	-0.710888	-0.737116	-0.73729
4	-0.869428	-0.877197	-0.970332	-0.904736	-0.890050	-0.966279	-0.90651
...
12718	-0.247231	-0.201970	-0.597255	-0.227293	-0.204002	-0.597073	-0.34921
12719	-1.018130	-0.788207	-0.989302	-0.918561	-0.844610	-0.985376	-0.86361
12720	0.394533	-0.466688	0.661086	-0.420848	-0.366561	0.669688	0.07291
12721	0.104957	0.865940	0.547267	1.390274	0.886261	0.523279	1.32421
12722	0.128436	-0.086534	0.616823	-0.144341	-0.124553	0.625129	-0.18341

12723 rows × 9 columns

In [47]:

trip[num_cols].describe()

Out[47]:

	start_scan_to_end_scan	actual_distance_to_destination	actual_time	osrm_time	osrm_distance	segment_actual_time_sum	segment_osrm_distai
count	1.272300e+04	1.272300e+04	1.272300e+04	1.272300e+04	1.272300e+04	1.272300e+04	1.272300e+04
mean	-1.808268e-17	-5.267735e-17	-6.830799e-17	7.396469e-17	5.696514e-17	-8.802900e-17	2.966514e-17
std	1.000039e+00	1.000039e+00	1.000039e+00	1.000039e+00	1.000039e+00	1.000039e+00	1.000039e+00
min	-1.162918e+00	-8.785574e-01	-1.065181e+00	-1.001514e+00	-9.229378e-01	-1.061764e+00	-9.375000e-01
25%	-7.207269e-01	-7.065920e-01	-7.363685e-01	-7.111809e-01	-7.077649e-01	-7.371165e-01	-7.228000e-01
50%	-3.411472e-01	-4.689012e-01	-4.012322e-01	-3.931975e-01	-4.836339e-01	-3.997380e-01	-4.628000e-01
75%	4.023595e-01	4.073375e-01	4.650634e-01	4.224989e-01	4.419548e-01	4.596223e-01	4.488000e-01
max	4.049455e+00	4.178358e+00	4.031419e+00	4.113871e+00	4.150641e+00	4.037107e+00	4.130000e+00

Recommendation

There is a significant difference between OSRM and actual parameters.

Hence,

- Revisit information fed to routing engine for trip planning. Check for discrepancies with transporters, if the routing engine is configured for optimum results.
- North, South and West Zones corridors have significant traffic of orders. But, we have a smaller presence in Central, Eastern and North-Eastern zone. However it would be difficult to conclude this, by looking at just 2 months data. It is worth investigating and increasing our presence in these regions.
- From state point of view, we have heavy traffic in Maharashtra followed by Karnataka. This is a good indicator that we need to plan for resources on ground in these 2 states on priority. Especially, during festive seasons.

In []: