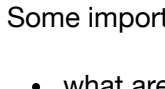


# Welcome to my EDA Kernel

## Description:

The dataset for this competition includes aggregate stopped vehicle information and intersection wait times. Your task is to predict congestion, based on an aggregate measure of stopping distance and waiting times, at intersections in 4 major US cities: Atlanta, Boston, Chicago & Philadelphia.



## Objective:

It's a first contact with the data, so I want to explore it and understand how the data is.

Some important things that is standard to analyze:

- what are the data types of the features?
- We have missing values?
- How many unique values we have in each feature;
- The shape of full dataset.
- The entropy of each feature (that show us the level of disorder on this column, it's like a "messy metric")

After this first analyze we can think in other questions to explore:

- Which distribution we have in our columns?
- Which are the most common cities?
- Which are the distribution of the stops, time, distances?
- How long is our date range?
- What are the distribution of the regions?

And many more questions;

**I'm near of grandmaster tier, so, if you find this kernel useful or interesting, please don't forget to upvote the kernel =)**

## Importing the Main Libraries to work with data

```
In [1]: import numpy as np # linear algebra
import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)
import osipy as sp
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

import plotly.graph_objs as go
import plotly.tools as tils
from plotly.offline import plot, init_notebook_mode
import cufflinks
import cufflinks as cf
import plotly.figure_factory as ff

from sklearn.decomposition import PCA
from sklearn.cluster import KMeans
from sklearn.metrics import mean_squared_error
from sklearn.ensemble import RandomForestRegressor
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import LabelEncoder
from sklearn.linear_model import LinearRegression

from functools import partial
from hyperopt import fmin, hp, tpe, Trials, space_eval, STATUS_OK, STATUS_RUNNING

import os
for dirname, _, filenames in os.walk('/kaggle/input'):
    for filename in filenames:
        print(os.path.join(dirname, filename))

# Any results you write to the current directory are saved as output.

/kaggle/input/bigquery-geotab-intersection-congestion/train.csv
/kaggle/input/bigquery-geotab-intersection-congestion/sample_submission.csv
/kaggle/input/bigquery-geotab-intersection-congestion/test.csv
/kaggle/input/bigquery-geotab-intersection-congestion/submission_metric_map.json
/kaggle/input/bigquery-geotab-intersection-congestion/BigQuery-Dataset-Access.md
```

## Importing datasets

```
In [2]: df_train = pd.read_csv('/kaggle/input/bigquery-geotab-intersection-congestion/train.csv')
df_test = pd.read_csv('/kaggle/input/bigquery-geotab-intersection-congestion/test.csv')
```

## Util functions

```
In [3]: def resumetable(df):
    print(f'Dataset Shape: {df.shape}')
    summary = pd.DataFrame(df.dtypes, columns=['dtypes'])
    summary = summary.reset_index()
    summary['Name'] = summary['index']
    summary = summary[['Name', 'dtypes']]
    summary['Missing'] = df.isnull().sum().values
    summary['Uniques'] = df.nunique().values
    summary['First Value'] = df.loc[0].values
    summary['Second Value'] = df.loc[1].values
    summary['Third Value'] = df.loc[2].values

    for name in summary['Name'].value_counts().index:
        summary.loc[summary['Name'] == name, 'Entropy'] = round(stats.entropy(df[name].value_counts(normalize=True), base=2), 2)

    return summary

def reduce_mem_usage(df, verbose=True):
    numerics = ['int16', 'int32', 'int64', 'float16', 'float32', 'float64']
    start_mem = df.memory_usage().sum() / 1024**2
    for col in df.columns:
        col_type = df[col].dtypes
        if col_type in numerics:
            c_min = df[col].min()
            c_max = df[col].max()
            if str(col_type)[:3] == 'int':
                if c_min > np.iinfo(np.int8).min and c_max < np.iinfo(np.int8).max:
                    df[col] = df[col].astype(np.int8)
                elif c_min > np.iinfo(np.int16).min and c_max < np.iinfo(np.int16).max:
                    df[col] = df[col].astype(np.int16)
                elif c_min > np.iinfo(np.int32).min and c_max < np.iinfo(np.int32).max:
                    df[col] = df[col].astype(np.int32)
                elif c_min > np.iinfo(np.int64).min and c_max < np.iinfo(np.int64).max:
                    df[col] = df[col].astype(np.int64)
            else:
                if c_min > np.iinfo(np.float16).min and c_max < np.iinfo(np.float16).max:
                    df[col] = df[col].astype(np.float16)
                elif c_min > np.iinfo(np.float32).min and c_max < np.iinfo(np.float32).max:
                    df[col] = df[col].astype(np.float32)
                else:
                    df[col] = df[col].astype(np.float64)
    end_mem = df.memory_usage().sum() / 1024**2
    if verbose: print('Mem. usage decreased to (%.2f) Mb ((%.1f) % reduction)' % (end_mem, 100 * (start_mem - end_mem) / start_mem))
    return df
```

## Summary of the data

```
In [4]: resumetable(df_train)

Dataset shape: (857409, 28)
```

	Name	dtypes	Missing	Uniques	First Value	Second Value	Third Value	Entropy
0	RowId	int64	0	857409	1920335	1920336	1920337	19.71
1	IntersectionId	int64	0	2539	0	0	0	10.63
2	Latitude	float64	0	4505	33.7917	33.7917	33.7917	11.41
3	Longitude	float64	0	4541	-84.43	-84.43	-84.43	11.43
4	EntryStreetName	object	8189	1707	Marietta Boulevard Northwest	Marietta Boulevard Northwest	Marietta Boulevard Northwest	8.97
5	ExitStreetName	object	5534	1693	Marietta Boulevard Northwest	Marietta Boulevard Northwest	Marietta Boulevard Northwest	8.95
6	EntryHeading	object	0	8	NW	SE	NW	2.85
7	ExitHeading	object	0	8	NW	SE	NW	2.85
8	Hour	int64	0	24	0	0	1	4.50
9	Weekend	int64	0	2	0	0	0	0.85
10	Month	int64	0	9	6	6	6	2.81
11	Path	object	0	15111	Marietta Boulevard Northwest, NW, Marietta Boule...	Marietta Boulevard Northwest, SE, Marietta Boule...	Marietta Boulevard Northwest, NW, Marietta Boule...	11.69
12	TotalTimeStopped_p20	int64	0	172	0	0	0	0.92
13	TotalTimeStopped_p40	int64	0	234	0	0	0	2.11
14	TotalTimeStopped_p50	int64	0	264	0	0	0	2.70
15	TotalTimeStopped_p60	int64	0	311	0	0	0	3.55
16	TotalTimeStopped_p80	int64	0	403	0	0	0	5.06
17	TimeFromFirstStop_p20	int64	0	239	0	0	0	0.98
18	TimeFromFirstStop_p40	int64	0	306	0	0	0	2.25
19	TimeFromFirstStop_p50	int64	0	329	0	0	0	2.87
20	TimeFromFirstStop_p60	int64	0	351	0	0	0	3.77
21	TimeFromFirstStop_p80	int64	0	355	0	0	0	5.35
22	DistanceToFirstStop_p20	float64	0	3479	0	0	0	1.35
23	DistanceToFirstStop_p40	float64	0	6257	0	0	0	3.23
24	DistanceToFirstStop_p50	float64	0	7483	0	0	0	4.16
25	DistanceToFirstStop_p60	float64	0	9495	0	0	0	5.55
26	DistanceToFirstStop_p80	float64	0	13267	0	0	0	8.10
27	City	object	0	4	Atlanta	Atlanta	Atlanta	1.85

Nice, this func give us a lot of cool and useful informations;

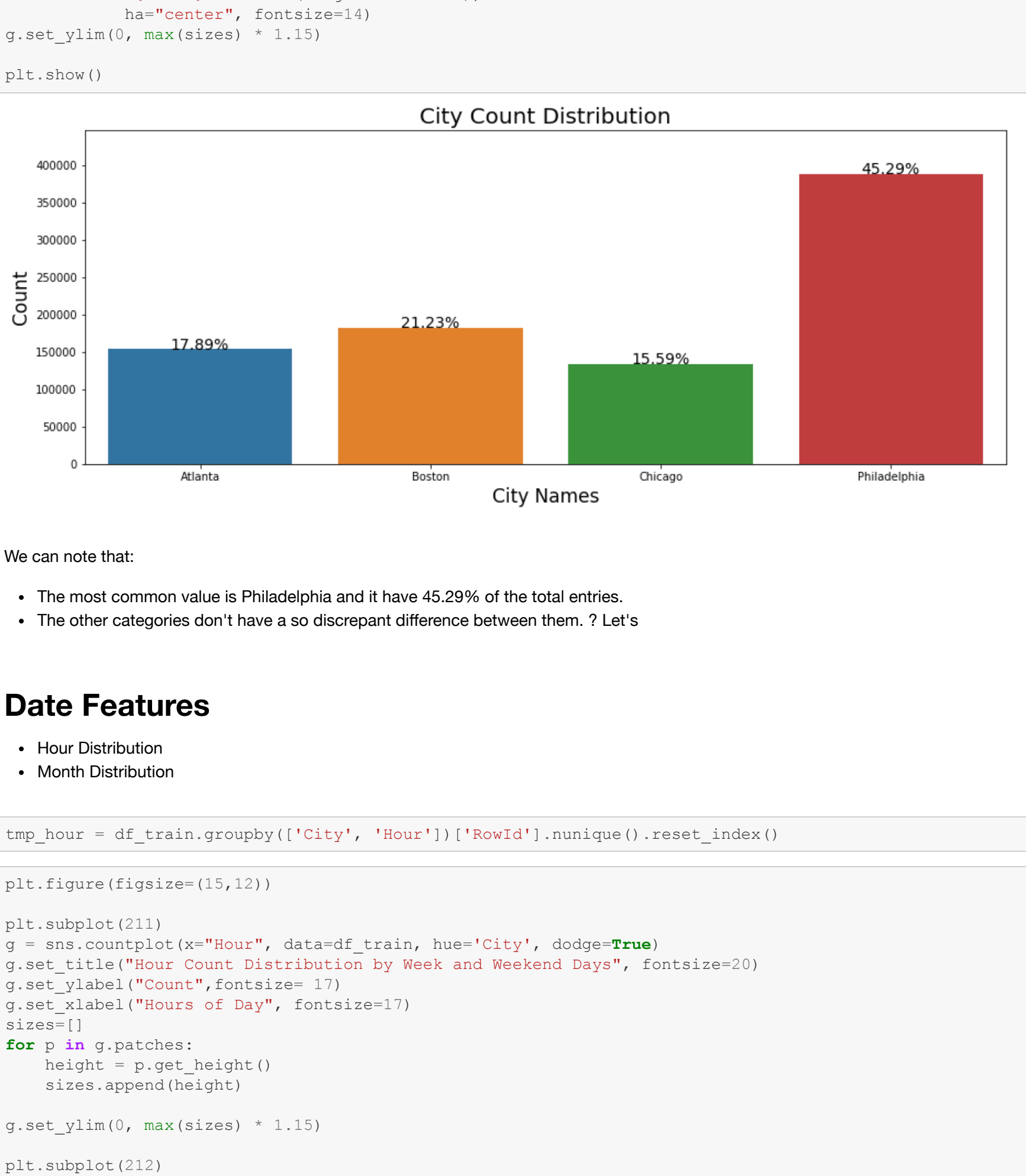
- We have only two features with missing values. Entry and Exit StreetName

## City's

- I will start exploring the distribution of City's because it is a categorical with only a few categories inside.

```
In [5]: resumetable(df_train)

Dataset shape: (857409, 28)
```



We can note that:

- The most common city is Philadelphia and it have 45.29% of the total entries.
- The other categories don't have a so discrepant difference between them. 7 Let's

## Date Features

- Hour Distribution
- Month Distribution

```
In [7]: tmp_hour = df_train.groupby(['City', 'Hour'])['RowId'].nunique().reset_index()
```

```
In [8]: plt.figure(figsize=(15,12))

plt.subplot(211)
g = sns.countplot(x="Hour", data=df_train, hue="City", dodge=True)
g.set_title("Hour Count Distribution by Week and Weekend Days", fontsize=20)
g.set_ylabel("Count", fontsize=17)
g.set_xlabel("Hours of Day", fontsize=17)
g.set_xticklabels(g.get_xticklabels(), rotation=45)
for p in g.patches:
    height = p.get_height()
    sizes.append(height)

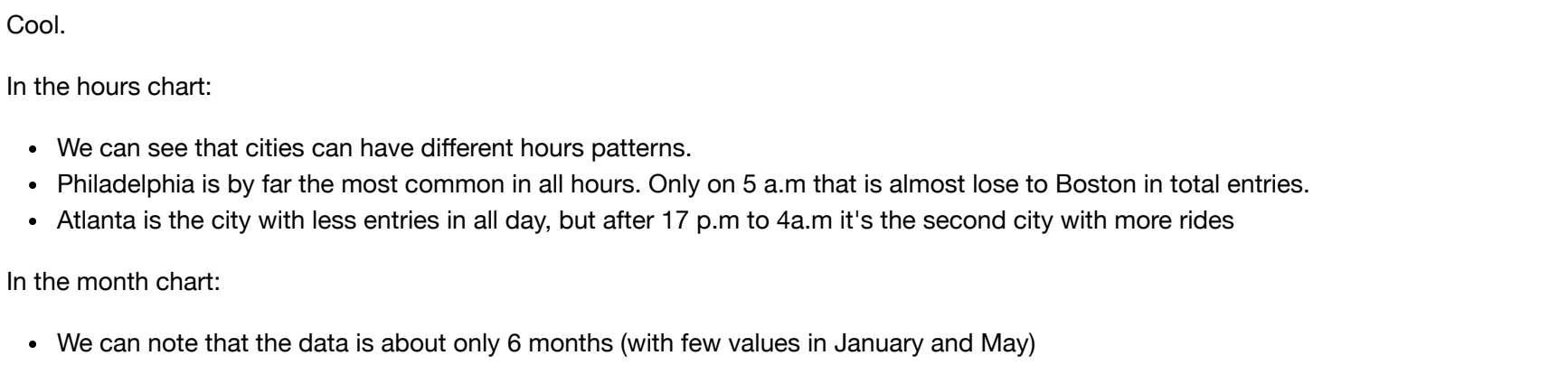
g.set_ylim(0, max(sizes) * 1.15)
```

```
plt.subplot(212)
g1 = sns.countplot(x="Month", data=df_train, hue="City", dodge=True)
g1.set_title("Hour Count Distribution by Week and Weekend Days", fontsize=20)
g1.set_ylabel("Count", fontsize=17)
g1.set_xlabel("Hours of Day", fontsize=17)
g1.set_xticklabels(g.get_xticklabels(), rotation=45)
for p in g1.patches:
    height = p.get_height()
    sizes.append(height)

g1.set_ylim(0, max(sizes) * 1.15)

plt.subplots_adjust(hspace = 0.3)

plt.show()
```



Cool.

- We can see that cities can have different hours patterns.
- Philadelphia is by far the most common in all hours. Only on 5 a.m that is almost lose to Boston in total entries.
- Atlanta is the city with less entries in all day, but after 17 p.m to 4 a.m it's the second city with more rides

In the hours chart:

- We can note that the data is about only 6 months (with few values in January and May)
- Also, the pattern of the Boston City improved throughout the time and the others seem very unchanged.

Now, let's explore the Entry and Exit features.

## EntryHeading and Exit Heading

```
In [9]: plt.figure(figsize=(15,12))

tmp = round((df_train.groupby(['EntryHeading'])['RowId'].nunique() / total) * 100).reset_index()

plt.subplot(211)
data=df_train,
order=list(tmp['EntryHeading'].values),
hue='ExitHeading', dodge=True)
g.set_title("Entry Heading by Exit Heading", fontsize=20)
g.set_ylabel("Count", fontsize=17)
g.set_xlabel("Entry Heading Region", fontsize=17)
gt = g.twinx()
gt = sns.pointplot(x="EntryHeading", y="RowId",
                    data=tmp, order=list(tmp['EntryHeading'].values),
                    color='black', legend=False)
gt.set_ylabel("% of Total (Black Line)", fontsize=16)
gt.set_xticklabels(gt.get_xticklabels(), rotation=45)
for p in g.patches:
    height = p.get_height()
    sizes.append(height)

g.set_ylim(0, max(sizes) * 1.15)
```

```
plt.subplot(212)
g1 = sns.countplot(x="EntryHeading", order=list(tmp['EntryHeading'].values),
                    data=df_train, hue="City",
                    fontsize=20)
g1.set_title("Entry Heading Distribution by Cities", fontsize=20)
g1.set_ylabel("Count", fontsize=17)
g1.set_xlabel("Entry Heading Region", fontsize=17)
for p in g1.patches:
    height = p.get_height()
    sizes.append(height)

g1.set_ylim(0, max(sizes) * 1.15)

plt.subplots_adjust(hspace = 0.3)

plt.show()
```



Nice.

In Entry and Exit Heading chart:

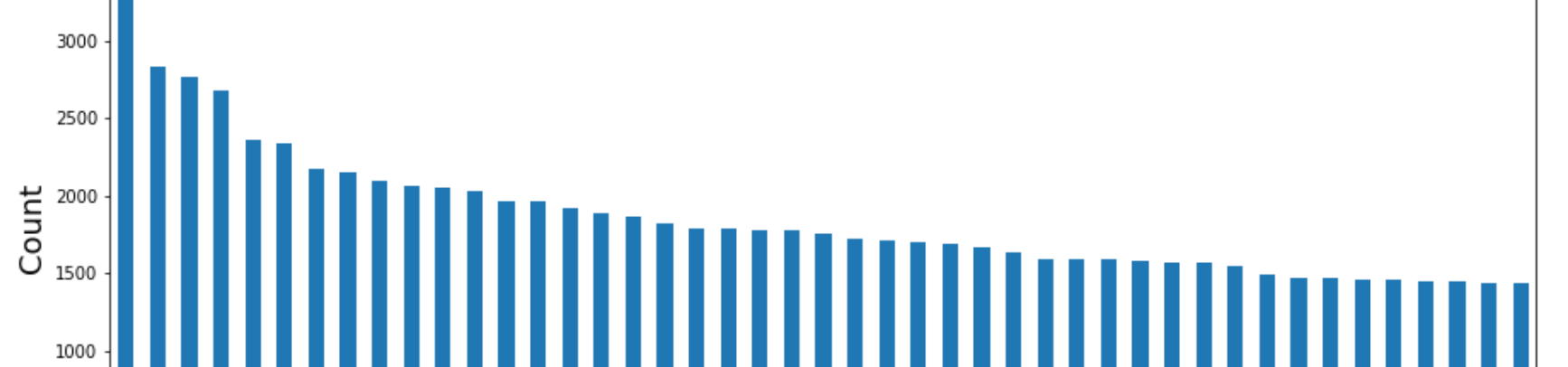
- We can note that in general the Entry and Exit Region is exactly the same.

In Entry by Cities chart:

- We can note the difference patterns on the cities. It's a very interesting and could give us many interesting insights.

## IntersectionID

```
In [10]: plt.figure(figsize=(15,6))
df_train.intersection_id.value_counts().plot(kind='bar')
plt.xlabel("Intersection Number", fontsize=18)
plt.ylabel("Count", fontsize=18)
plt.title("TOP 45 most common IntersectionID's ", fontsize=22)
plt.show()
```



```
In [11]: df_train.groupby(['IntersectionId', 'EntryHeading', 'ExitHeading'])['RowId'].count().reset_index().head()
```

```
Out [11]:
```

IntersectionId	EntryHeading	ExitHeading	RowId	
0	0	E	150	
1	0	E	N	24
2	0	N	E	31
3	0	N	N	59
4	0	N	W	22

## Exploring numerical features

The targets are:

- TotalTimeStopped\_p20
- TotalTimeStopped\_p50
- TotalTimeStopped\_p80
- DistanceToFirstStop\_p20
- DistanceToFirstStop\_p50
- DistanceToFirstStop\_p80

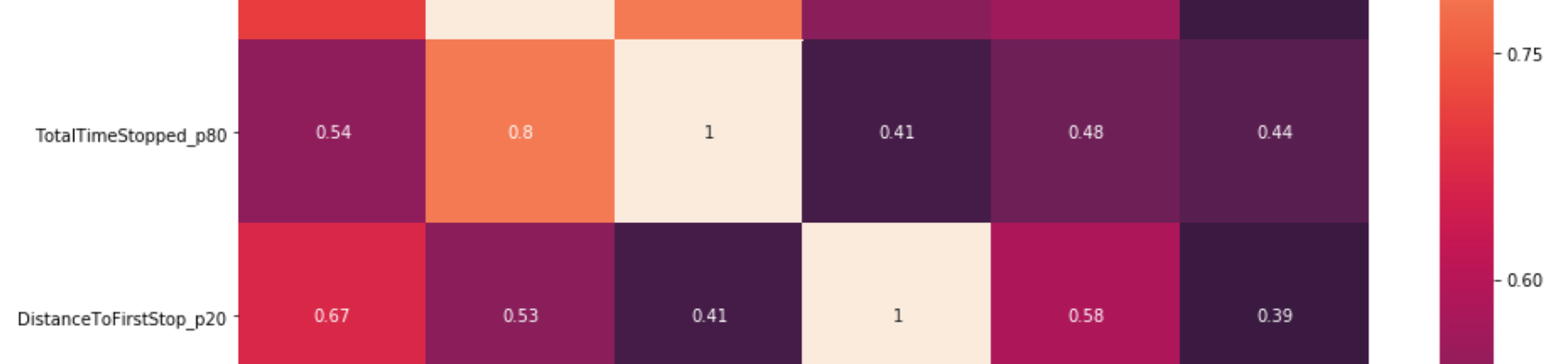
And the as the TimeFromFirstStop is an optional data, I will use it to see the correlations.

```
In [12]: t_stopped = ['TotalTimeStopped_p20',
                    'TotalTimeStopped_p50',
                    'TotalTimeStopped_p80']
t_first_stopped = ['TimeFromFirstStop_p20',
                   'TimeFromFirstStop_p50',
                   'TimeFromFirstStop_p80']
d_first_stopped = ['DistanceToFirstStop_p20',
                   'DistanceToFirstStop_p50',
                   'DistanceToFirstStop_p80']
```

## Heatmap Target Features

```
In [13]: plt.figure(figsize=(15,12))
plt.title("Correlation of Features for Train Set", fontsize=22)
sns.heatmap(df_train[t_stopped +
                    t_first_stopped +
                    d_first_stopped].astype(float).corr(),
            vmax=1.0, annot=True)

plt.show()
```



Cool!

We can see that the best correlation between the metrics are:

- Distance to First Stop p20 and Total Time Stopped p20 have a high correlation.

## Scaling the target

- Getting the min\_max transformation to get clusterization and PCA features

```
In [14]: from sklearn.preprocessing import minmax_scale

target_cols = t_stopped + d_first_stopped
```

```
In [15]: for col in target_cols:
    df_train[col+str("minmax")] = minmax_scale(df_train[col], feature_range=(0,1))

min_max_cols = ['TotalTimeStopped_p20_minmax', 'TotalTimeStopped_p50_minmax',
                'TotalTimeStopped_p80_minmax', 'DistanceToFirstStop_p20_minmax',
                'DistanceToFirstStop_p50_minmax', 'DistanceToFirstStop_p80_minmax']
```

## PCA

- To better see the distribution of our metrics, lets apply PCA to reduce the dimensionality of the data

```
In [16]: pca = PCA(n_components=3, random_state=5)
principalComponents = pca.fit_transform(df_train[min_max_cols])
principalDf = pd.DataFrame(principalComponents)
# df.drop(cols, axis=1, inplace=True)
prefix="target_PCA"
principalDf.rename(columns=lambda x: str(prefix)+str(x), inplace=True)
df_train = pd.concat([df_train, principalDf], axis=1)
```

Nice, now we have the PCA features... Let's see the ratio of explanation of the first two Principal Components

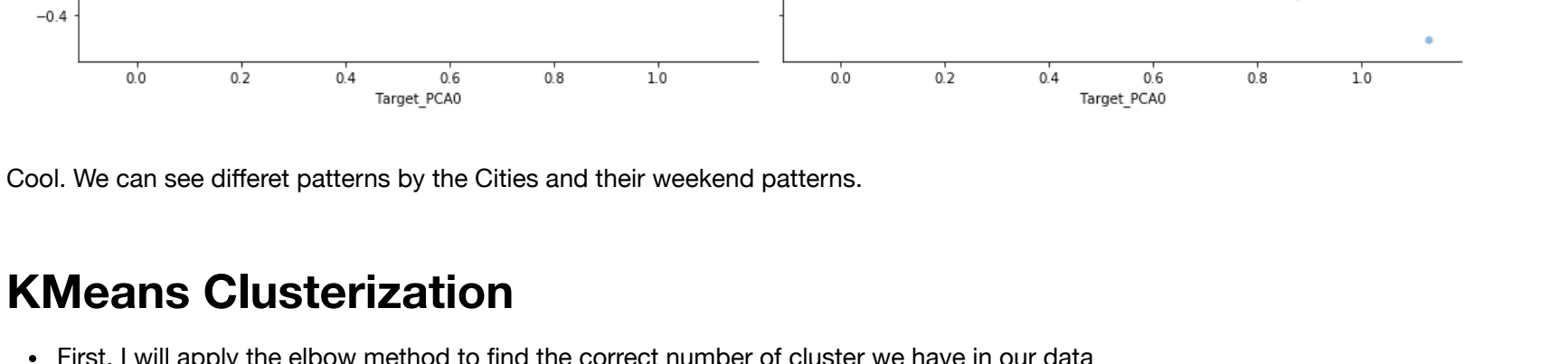
```
In [17]: pca.explained_variance_ratio_[0:2].sum()
```

```
Out [17]: 0.8393288157969666
```

With the 2 first components we have almost 84% of the data explained. It's a very way to easiest visualize the differences between the patterns.

## Scatter plot of cities by the PCA

```
In [18]: g = sns.FacetGrid(df_train.sample(50000), col="City",
                        col_wrap=4, height=5, aspect=1.5, hue="Weekend")
g.map(sns.scatterplot, "target_PCA0", "target_PCA1", alpha=.5).add_legend()
g.set_titles('toot_names', fontsize=17)
plt.show()
```



Cool. We can see different patterns by the Cities and their weekend patterns.

## KMeans Clusterization

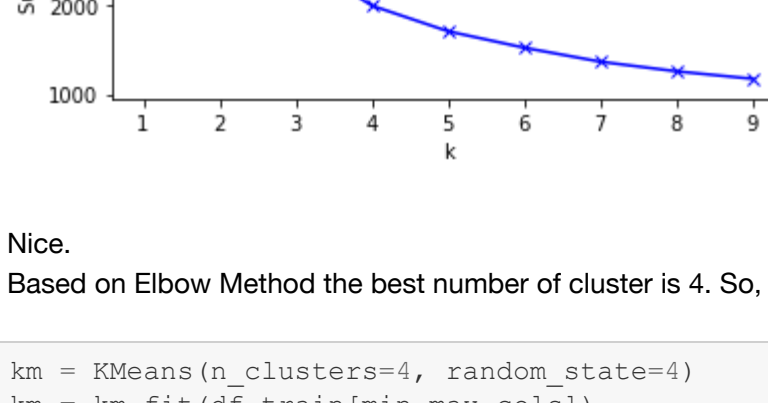
- First, I will apply the elbow method to find the correct number of cluster we have in our data
- After it, we will implement the kmeans with the best quantity

```
In [19]: #sum of squared distances
ssd=[]

K = range(1,10)

for k in K:
    km = KMeans(n_clusters=k, random_state=4)
    km = km.fit(df_train[min_max_cols])
    ssd.append(km.inertia_)

plt.plot(K, ssd, 'bx-')
plt.xlabel('k')
plt.ylabel('Sum of squared distances')
plt.title('Elbow Method For Optimal k')
plt.show()
```



Nice.

Based on Elbow Method the best number of cluster is 4. So, let's apply the K means on data.

```
In [20]: km = KMeans(n_clusters=4, random_state=4)
km = km.fit(df_train[min_max_cols])
df_train['clusters_k'] = km.predict(df_train[min_max_cols])
```

## Plotting Clusters

- Understanding the cluster distribution
- Exploring by Cities



```
In [21]: tmp = pd.crosstab([train['City'], df_train['clusters_T']],
                        df_train['clusters_T']).reset_index().rename(columns={0:'perc'})

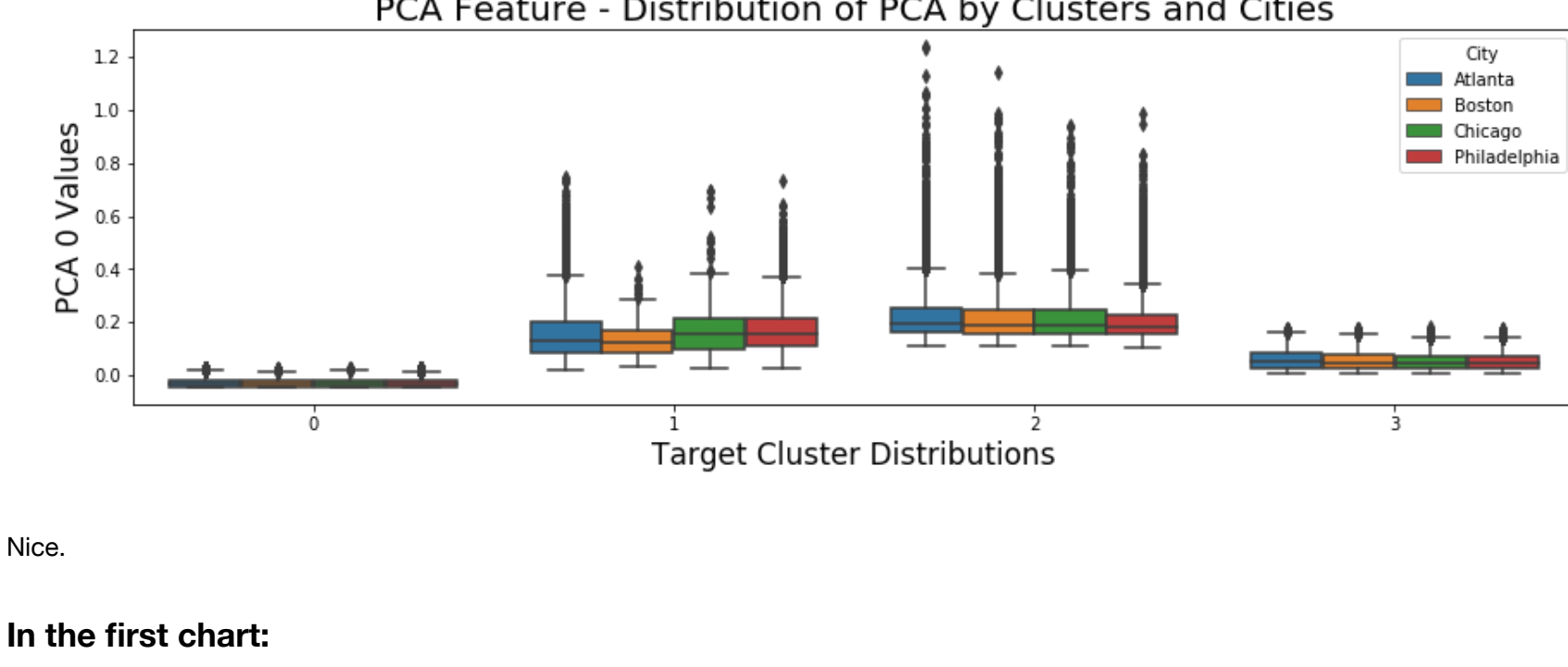
total = len(df_train)
plt.figure(figsize=(15,16))

plt.subplot(311)
g = sns.countplot(x='clusters_T', data=df_train)
g.set_title("Cluster Target Count Distribution", fontsize=20)
g.set_ylabel("Count", fontsize=17)
g.set_xlabel("Target Cluster Distributions", fontsize=17)
sizes=[]
for p in g.patches:
    height = p.get_height()
    sizes.append(height)
    g.text(p.get_x()+p.get_width()/2.,
          height + 3,
          '{1.2f}'.format(height/total*100),
          ha='center', fontsize=14)
g.set_ylim(0, max(sizes) * 1.15)

plt.subplot(312)
g1 = sns.countplot(x='clusters_T', data=df_train, hue='City')
g1.set_title("CITIES - Cluster Target Distribution", fontsize=20)
g1.set_ylabel("Count", fontsize=17)
g1.set_xlabel("Target Cluster Distributions", fontsize=17)
sizes=[]
for p in g1.patches:
    height = p.get_height()
    sizes.append(height)
    g1.text(p.get_x()+p.get_width()/2.,
           height + 3,
           '{1.2f}'.format(height/total*100),
           ha='center', fontsize=10)
g1.set_ylim(0, max(sizes) * 1.15)

plt.subplot(313)
g1 = sns.boxplot(x='clusters_T', y='Target_PCA0',
                data=df_train, hue='City')
g1.set_title("PCA Feature - Distribution of PCA by Clusters and Cities",
            fontsize=20)
g1.set_ylabel("PCA 0 Values", fontsize=17)
g1.set_xlabel("Target Cluster Distributions", fontsize=17)
plt.subplots_adjust(hspace = 0.5)

plt.show()
```



Nice.

**In the first chart:**

- We can note that the most common cluster is the 1 that have 73% of all data.

**Second chart:**

- Philadelphia is the most common in the first 3 clusters.
- Boston is the second most common in 0,1 and the most common on Cluster 3;
- In the second cluster, Atlanta is the second most common city.

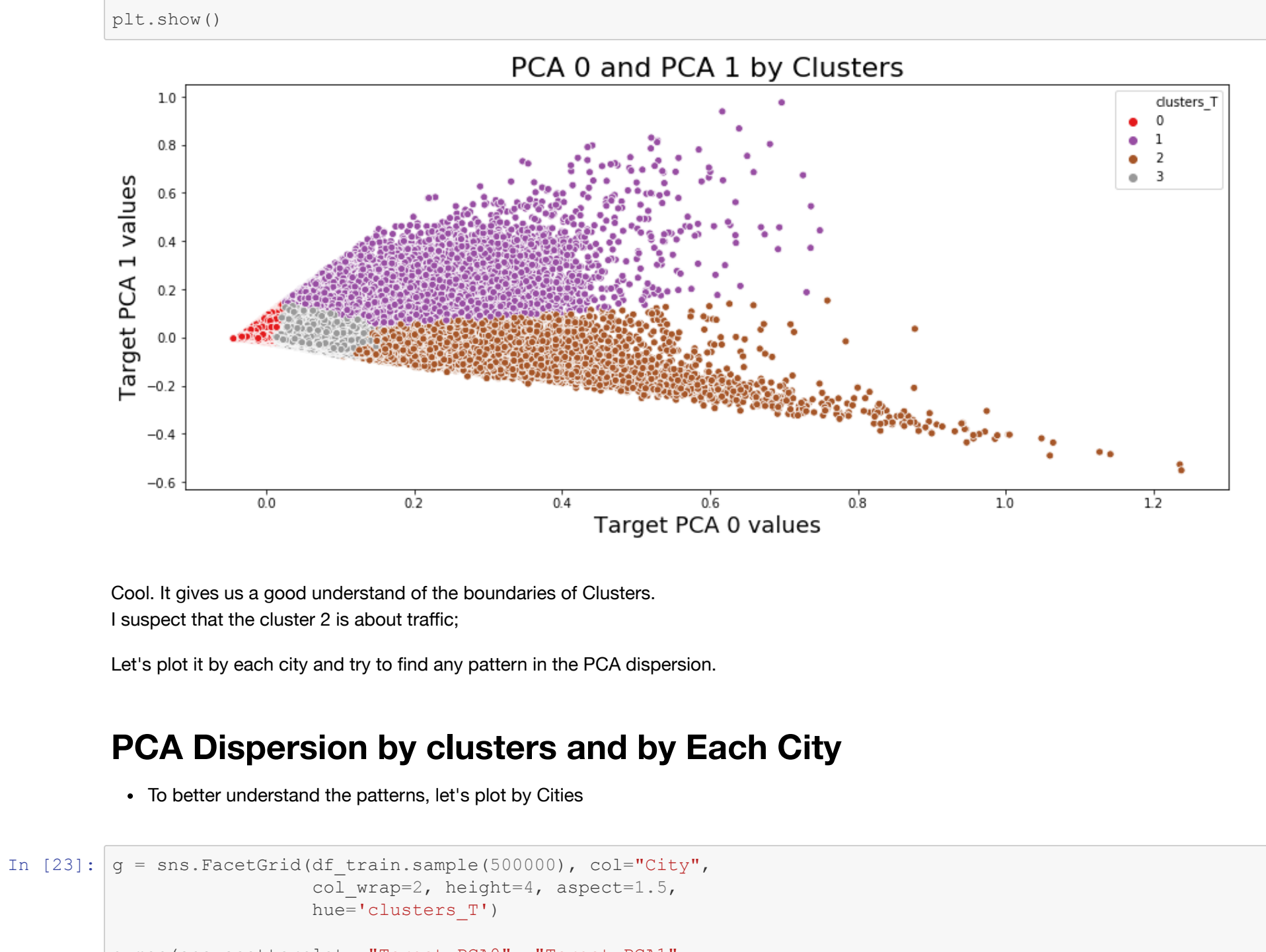
**Third Chart:**

- Is clear to understand how the algorithm divided the data in PCA values

**NOTE: EVERY TIME I RUN IT, THE VALUES CHANGES, SO SORRY BY THE WRONG**

## PCA values by CLUSTERS

- Let's see in another way how the algorithm divided the data by the clusterization

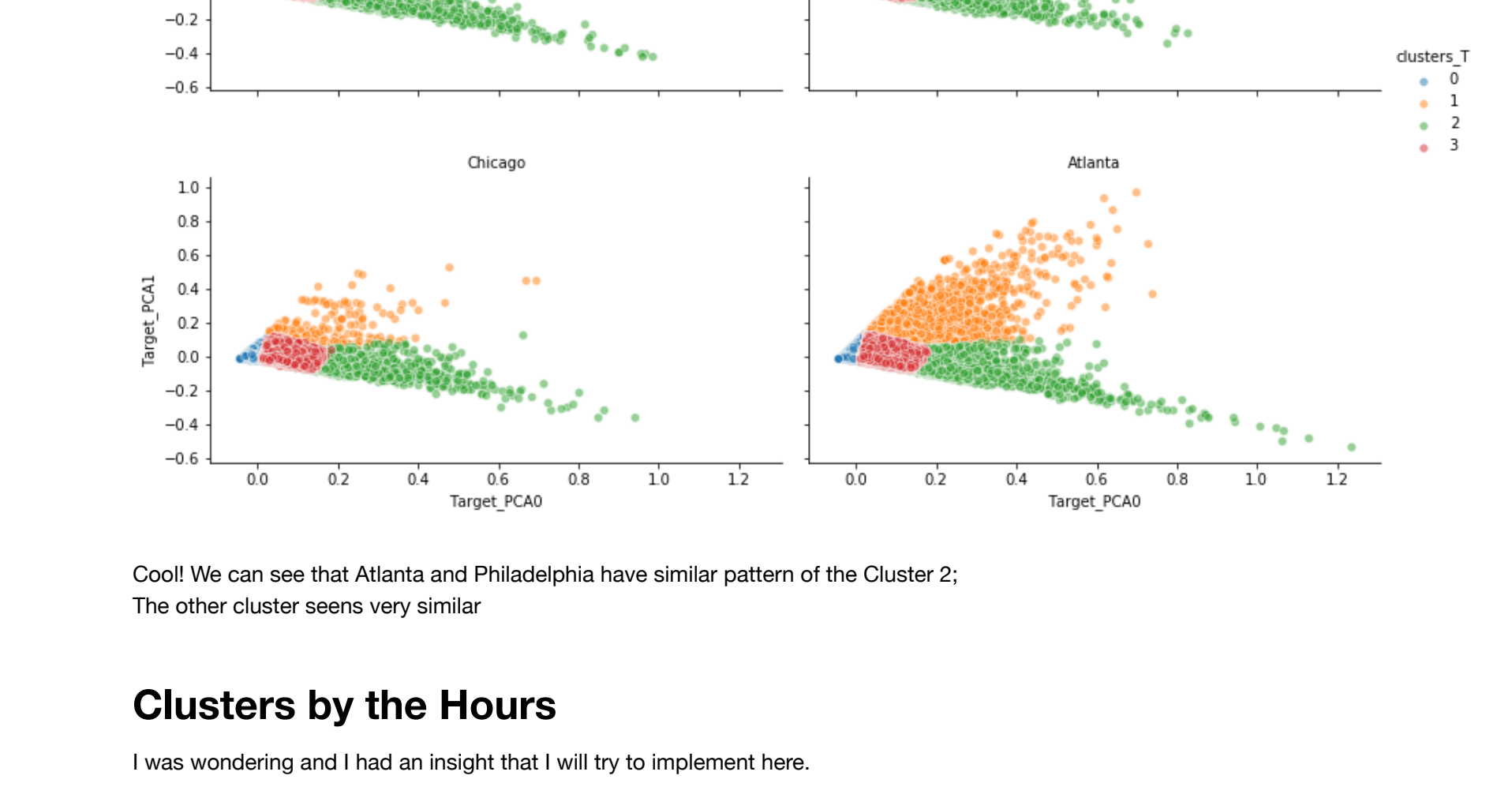


Cool. It gives us a good understanding of the boundaries of Clusters.  
I suspect that the cluster 2 is about traffic.

Let's plot it by each city and try to find any pattern in the PCA dispersion.

## PCA Dispersion by clusters and by Each City

- To better understand the patterns, let's plot by Cities

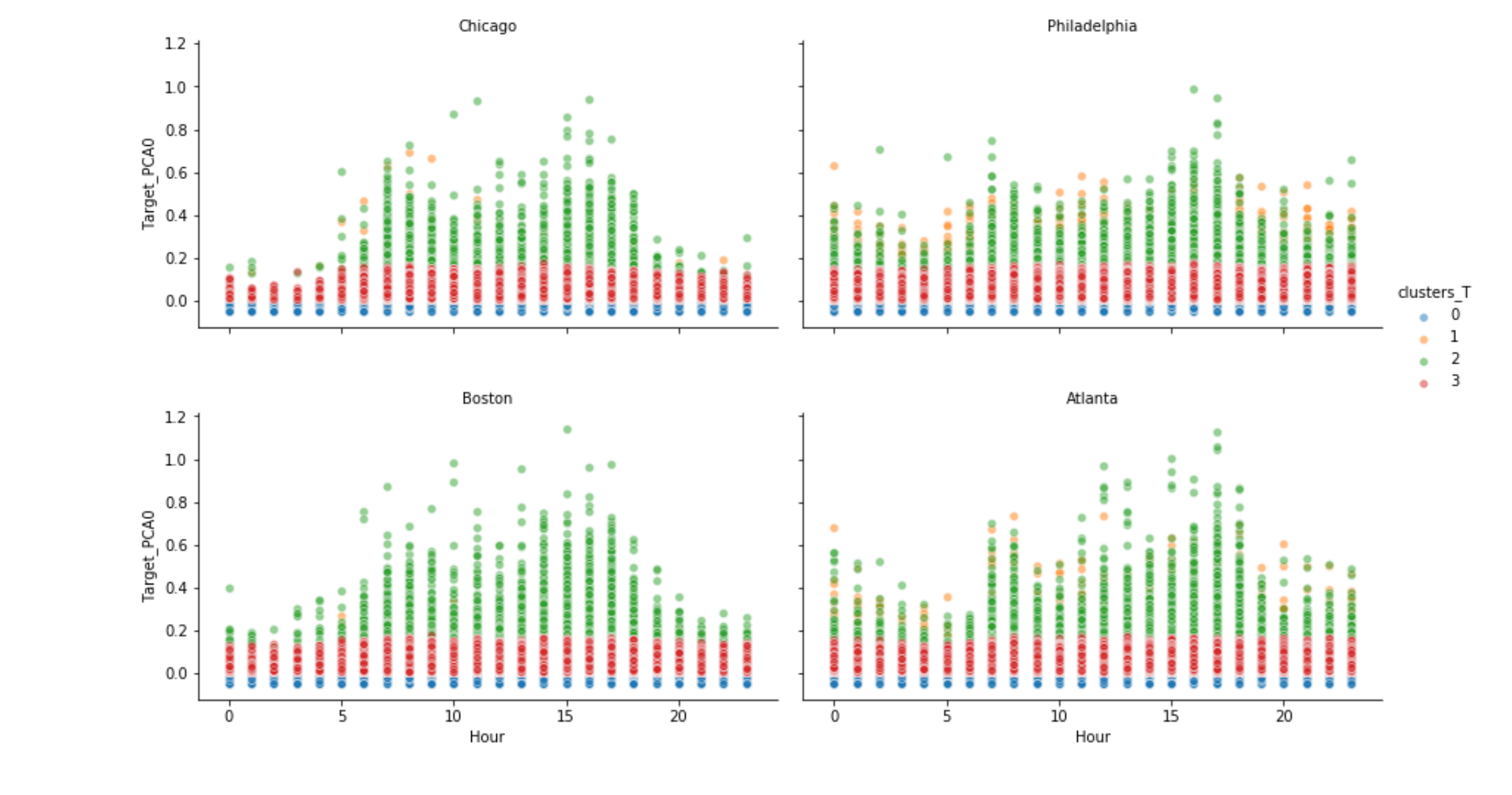


Cool! We can see that Atlanta and Philadelphia have similar pattern of the Cluster 2.  
The other cluster seems very similar

## Clusters by the Hours

I was wondering and I had an insight that I will try to implement here.

- I think that make a lot of sense explore the hours by the clusters
- Let's see the distribution of PCA0 and the Clusters by the Hours



Cool! We can have a best intuition about the data and how it possible clustered the data.

```
In [25]: round(pd.crosstab([df_train['clusters_T'], df_train['Weekend']],
                        df_train['clusters_T']).reset_index().rename(columns={0:'perc'}),
          2)
```

```
Out [25]:
```

	City	Atlanta	Boston	Chicago	Philadelphia
clusters_T	Weekend				
	0	0	15.0	22.0	22.0
0	1	22.0	16.0	3.0	58.0
	0	24.0	2.0	3.0	71.0
1	1	30.0	1.0	0.0	69.0
	0	24.0	32.0	18.0	28.0
2	1	43.0	15.0	2.0	41.0
	0	17.0	26.0	19.0	38.0
3	1	27.0	14.0	2.0	56.0

## Modeling

- As I was getting problems with my model, I decided to implement the solution of the public kernels
- I will import the datasets again

Many parts of this implementation I got on [@dcaichara](#) Kernel.  
You can see the kernel here: <https://www.kaggle.com/dcaichara/feature-engineering-and-lightgbm>

```
In [26]: df_train = pd.read_csv('/kaggle/input/biqquery-geotab-intersection-congestion/train.csv')
df_test = pd.read_csv('/kaggle/input/biqquery-geotab-intersection-congestion/test.csv')
```

## Hour Feature

- Let's encode the Hour Features

```
In [27]: def date_cyc_enc(df, col, max_val):
df[col + "_sin"] = np.sin(2 * np.pi * df[col]/max_val)
df[col + "_cos"] = np.cos(2 * np.pi * df[col]/max_val)
return df

df_train = date_cyc_enc(df_train, 'Hour', 24)
df_test = date_cyc_enc(df_test, 'Hour', 24)
```

## Flag - is day?

Testing some features about the data

```
In [28]: df_train['is_day'] = df_train['Hour'].apply(lambda x: 1 if 7 < x < 18 else 0)
df_test['is_day'] = df_test['Hour'].apply(lambda x: 1 if 7 < x < 18 else 0)

df_train['is_morning'] = df_train['Hour'].apply(lambda x: 1 if 6 < x < 10 else 0)
df_test['is_morning'] = df_test['Hour'].apply(lambda x: 1 if 6 < x < 10 else 0)

df_train['is_night'] = df_train['Hour'].apply(lambda x: 1 if 17 < x < 20 else 0)
df_test['is_night'] = df_test['Hour'].apply(lambda x: 1 if 17 < x < 20 else 0)

df_train['is_day_weekend'] = np.where((df_train['is_day'] == 1) & (df_train['Weekend'] == 1), 1, 0)
df_test['is_day_weekend'] = np.where((df_test['is_day'] == 1) & (df_test['Weekend'] == 1), 1, 0)

df_train['is_mor_weekend'] = np.where((df_train['is_morning'] == 1) & (df_train['Weekend'] == 1), 1, 0)
df_test['is_mor_weekend'] = np.where((df_test['is_morning'] == 1) & (df_test['Weekend'] == 1), 1, 0)

df_train['is_nig_weekend'] = np.where((df_train['is_night'] == 1) & (df_train['Weekend'] == 1), 1, 0)
df_test['is_nig_weekend'] = np.where((df_test['is_night'] == 1) & (df_test['Weekend'] == 1), 1, 0)
```

## Intersec - concatenating IntersectionId and City

```
In [29]: df_train['Intersec'] = df_train['IntersectionId'].astype(str) + df_train['City']
df_test['Intersec'] = df_test['IntersectionId'].astype(str) + df_test['City']

print(df_train["Intersec"].sample(6).values)

['1786Philadelphia' '1525Philadelphia' '1823Chicago'
 '1786Philadelphia' '175Atlanta']
```

## Label Encoder of Intersec + City

```
In [30]: le = LabelEncoder()

le.fit(pd.concat([df_train["Intersec"], df_test["Intersec"]]).drop_duplicates().values)
df_train["Intersec"] = le.transform(df_train["Intersec"])
df_test["Intersec"] = le.transform(df_test["Intersec"])
```

## Street Feature

- Extracting informations from street features

```
In [31]: road_encoding = {
    'Road': 1,
    'Street': 2,
    'Avenue': 2,
    'Drive': 3,
    'Boulevard': 4,
    'Boulevard': 4
}

def encode(x):
    if pd.isna(x):
        return 0
    for road in road_encoding.keys():
        if road in x:
            return road_encoding[road]
    return 0
```

## Creating the new feature

```
In [33]: df_train['EntryType'] = df_train['EntryStreetName'].apply(encode)
df_train['ExitType'] = df_train['ExitStreetName'].apply(encode)
df_test['EntryType'] = df_test['EntryStreetName'].apply(encode)
df_test['ExitType'] = df_test['ExitStreetName'].apply(encode)
```

## Encoding the Regions

```
In [34]: directions = {
    'N': 0,
    'NE': 1/4,
    'E': 1/2,
    'SE': 3/4,
    'S': 1,
    'SW': 3/4,
    'W': 1/2,
    'NW': 1/4
}
```

## Applying the transformation in Entry and Exit Heading Columns

```
In [35]: df_train['EntryHeading'] = df_train['EntryHeading'].map(directions)
df_train['ExitHeading'] = df_train['ExitHeading'].map(directions)

df_test['EntryHeading'] = df_test['EntryHeading'].map(directions)
df_test['ExitHeading'] = df_test['ExitHeading'].map(directions)
```

## Difference between the regions

```
In [36]: df_train['diffHeading'] = df_train['EntryHeading']-df_train['ExitHeading']
df_test['diffHeading'] = df_test['EntryHeading']-df_test['ExitHeading']
```

## Getting the binary if the entry and exit was in the same street

```
In [37]: df_train['same_str'] = (df_train['EntryStreetName'] == df_train['ExitStreetName']).astype(int)
df_test['same_str'] = (df_test['EntryStreetName'] == df_test['ExitStreetName']).astype(int)
```

## Concatenating City and Month

```
In [38]: # Concatenating the city and month into one variable
df_train['city_month'] = df_train['City'] + df_train['Month'].astype(str)
df_test['city_month'] = df_test['City'] + df_test['Month'].astype(str)
```

## Month rainfall ratio by city and seasons

```
In [39]: monthly_rainfall = {'Atlanta': 5.02, 'Atlanta5': 3.95, 'Atlanta6': 3.63, 'Atlanta7': 5.12,
                          'Atlanta8': 3.67, 'Atlanta9': 4.09, 'Atlanta10': 3.11, 'Atlanta11': 4.10,
                          'Atlanta12': 3.82, 'Boston1': 3.82, 'Boston5': 3.30, 'Boston9': 3.24, 'Boston10': 3.22,
                          'Boston7': 3.06, 'Boston8': 3.37, 'Boston9': 3.47, 'Boston10': 3.79,
                          'Boston12': 3.38, 'Boston11': 3.73, 'Chicago1': 1.75, 'Chicago5': 3.38,
                          'Chicago6': 3.63, 'Chicago7': 3.51, 'Chicago8': 4.62, 'Chicago9': 3.27,
                          'Chicago10': 2.71, 'Chicago11': 3.01, 'Chicago12': 2.43,
                          'Philadelphia': 3.32, 'Philadelphia5': 3.88, 'Philadelphia6': 3.28,
                          'Philadelphia7': 4.39, 'Philadelphia8': 3.82, 'Philadelphia9': 3.89,
                          'Philadelphia10': 2.75, 'Philadelphia11': 3.16, 'Philadelphia12': 3.31}

# Creating a new column by mapping the city_month variable to it's corresponding average monthly rainfall
df_train['average_rainfall'] = df_train['city_month'].map(monthly_rainfall)
df_test['average_rainfall'] = df_test['city_month'].map(monthly_rainfall)
```

## Getting Dummies

```
In [40]: print(f'Shape before dummy transformation: {df_train.shape}')
df_train = pd.get_dummies(df_train, columns=['City'],
                        prefix='City', drop_first=False)
print(f'Shape after dummy transformation: {df_train.shape}')

df_test = pd.get_dummies(df_test, columns=['City'],
                        prefix='City', drop_first=False)
print(f'Shape after dummy transformation: {df_test.shape}')
Shape before dummy transformation: (857409, 43)
Shape after dummy transformation: (857409, 46)
```

## MinMax Scaling the lat and long

```
In [41]: from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
for col in ['Latitude', 'Longitude']:
    scaler.fit(df_train[col].values.reshape(-1, 1))
df_train[col] = scaler.transform(df_train[col].values.reshape(-1, 1))
df_test[col] = scaler.transform(df_test[col].values.reshape(-1, 1))
```

## Dropping not used features

```
In [42]: df_train.drop(['RowId', 'Path', 'EntryStreetName', 'ExitStreetName',
                    'axis=1, inplace=True)
df_test.drop(['RowId', 'Path',
              'EntryStreetName', 'ExitStreetName'], axis=1, inplace=True)
```

```
In [43]: interesting_feat = ['IntersectionId', 'Latitude', 'Longitude', 'EntryHeading',
                          'ExitHeading', 'Hour', 'Weekend', 'Month',
                          'is_morning', 'is_night', 'is_day_weekend', 'is_mor_weekend',
                          'is_nig_weekend', 'Hour_sin',
                          'Hour', 'same_str', 'Intersec', 'EntryType',
                          'ExitType', 'diffHeading', 'average_rainfall', 'is_day',
                          'City_Boston', 'City_Chicago', 'City_Philadelphia',
                          'City_Atlanta']

total_time = ['TotalTimeStopped_p20',
              'TotalTimeStopped_p50',
              'TotalTimeStopped_p80']

target_stopped = ['DistanceToFirstStop_p20',
                  'DistanceToFirstStop_p50',
                  'DistanceToFirstStop_p80']
```

## Setting X and y

```
In [44]: X = df_train[interesting_feat]
y = df_train[target_stopped]

X_test = df_test[interesting_feat]
```

```
In [45]: print(f'Shape of X: {X.shape}')
print(f'Shape of X_test: {X_test.shape}')

Shape of X: (857409, 25)
Shape of X_test: (192035, 25)
```

## Reduce memory usage

X = reduce\_mem\_usage(X)  
X\_test = reduce\_mem\_usage(X\_test)

## Splitting data into train and validation

```
In [46]: X_train, X_val, y_train, y_val = train_test_split(X, y, test_size=0.10,
                                                    random_state=42)
```

## Hyperopt Space

- Here we will set all range of our hyperparameters

```
In [47]: # Define searched space
hyper_space = {'objective': 'regression',
               'metric': 'mae',
               'boosting': 'gbdt',
               'n_estimators': hp.choice('n_estimators', [25, 40, 50, 75, 100, 250, 500]),
               'max_depth': hp.choice('max_depth', list(range(6, 18, 2))),
               'num_leaves': hp.choice('num_leaves', list(range(20, 180, 20))),
               'subsample': hp.choice('subsample', [0.7, 0.8, 0.9, 1]),
               'colsample_bytree': hp.uniform('colsample_bytree', 0.7, 1),
               'learning_rate': hp.uniform('learning_rate', 0.03, 0.12),
               'reg_alpha': hp.choice('reg_alpha', [1, 2, 3, 4, 5, 6]),
               'reg_lambda': hp.choice('reg_lambda', [1, 1.2, 1.3, 1.4, 1.5, 1.6]),
               'min_child_samples': hp.choice('min_child_samples', [20, 45, 70, 100])}
```

## Building Hyperopt Function to be optimized

```
In [48]: cat_feat = ['IntersectionId', 'Hour', 'Weekend', 'Month',
                  'is_day', 'is_morning', 'is_night',
                  'same_str', 'Intersec', 'City_Atlanta', 'City_Boston',
                  'City_Chicago', 'City_Philadelphia', 'EntryType', 'ExitType']
```

```
In [49]: from sklearn.model_selection import KFold
import lightgbm as lgb

def evaluate_metrics(params):

    all_preds_test = [0] * len(preds_test)

    print(f'Params: {params}')
    FOLDS = 4

    count = 1

    for i in range(len(all_preds_test)):

        score_mean = 0

        kf = KFold(n_splits=FOLDS, shuffle=False,
                  random_state=42)

        for tr_idx, val_idx in kf.split(X, y):

            X_tr, y_tr = X.iloc[tr_idx:], y.iloc[val_idx:]
            X_val, y_val = X.iloc[tr_idx:], y.iloc[val_idx:]

            lgb_train = lgb.Dataset(X_tr, label=y_tr.iloc[:,1])
            lgb_val = lgb.Dataset(X_val, label=y_val.iloc[:,1])

            lgbm_reg = lgb.train(params, lgb_train, 2000, valid_sets = [lgb_val],
                                categorical_feature=cat_feat,
                                verbose_eval=0,
                                early_stopping_rounds = 300)

            pred_lgb = lgbm_reg.predict(X_val, num_iteration=lgbm_reg.best_iteration)
            all_preds_test[i] = pred_lgb

        score_mean = np.sqrt(mean_squared_error(pred[0].values, y_val_sc[0].values))
        #score = metric(df_val, pred)

        print(f'Score Validation : {score.uniq}')

    pred = pd.DataFrame(all_preds_test).stack()
    pred = pd.DataFrame(pred)

    y_val_sc = pd.DataFrame(y_val).stack()
    y_val_sc = pd.DataFrame(y_val_sc)

    count = count + 1

    score = np.sqrt(mean_squared_error(pred[0].values, y_val_sc[0].values))
    #score = metric(df_val, pred)

    print(f'Full Score Run: {score}')

    return {
        'loss': score,
        'status': STATUS_OK
    }
```

## Running the hyperopt Function



```
MX_EVALS=15  
fit_tree_param_estimate  
best_vals = fmin(evaluate_metric,  
                space=hyper_space,  
                objectives=[],  
                algo=tpe.suggest,  
                max_evals=MAY_EVALS)  
  
# Print best parameters  
best_params = space_ova(hyper_space, best_vals)  
  
Params: {'boosting': 'gbdt', 'coarsample_bytree': 0.924765974717267, 'gpu_device_id': 0, 'learning_rate': 0.045508240433874, 'max_depth': 12, 'metric': 'rmse', 'mini_child_samples': 45, 'num_leaves': 1  
47}, 'objective': 'regression', 'subsample': 0.8}  
Oblivious | 0 / 15 [0:00:00<] 217.8s, best loss: ?]  
  
/opt/conda/lib/python3.6/site-packages/lightgbm/basic.py:1209: UserWarning:  
categorical_feature in Dataset is overridden.  
New categorical feature is ['City_Atlanta', 'City_Boston', 'City_Chicago', 'City_Philadelphia', 'Entr  
yType', 'ExitType', 'Hour', 'Intersec', 'IntersectionId', 'Month', 'Weekend', 'is_day', 'is_morning'  
, 'is_night', 'same_str']  
  
/opt/conda/lib/python3.6/site-packages/lightgbm/basic.py:762: UserWarning:  
categorical_feature in param dict is overridden.  
  
Score Validation : 6.40195761384795  
0bl | 0 / 15 [0:02:16<] 717.8s, best loss: ?]  
  
/opt/conda/lib/python3.6/site-packages/lightgbm/basic.py:1209: UserWarning:  
categorical_feature in Dataset is overridden.  
New categorical feature is ['City_Atlanta', 'City_Boston', 'City_Chicago', 'City_Philadelphia', 'Ent  
ryType', 'ExitType', 'Hour', 'Intersec', 'IntersectionId', 'Month', 'Weekend', 'is_day', 'is_morning'  
, 'is_night', 'same_str']  
  
/opt/conda/lib/python3.6/site-packages/lightgbm/basic.py:762: UserWarning:  
categorical_feature in param dict is overridden.  
  
Score Validation : 12.646824627955  
0bl | 0 / 15 [0:13:35<] 717.8s, best loss: ?]  
  
/opt/conda/lib/python3.6/site-packages/lightgbm/basic.py:1209: UserWarning:  
categorical_feature in Dataset is overridden.  
New categorical feature is ['City_Atlanta', 'City_Boston', 'City_Chicago', 'City_Philadelphia', 'Ent  
ryType', 'ExitType', 'Hour', 'Intersec', 'IntersectionId', 'Month', 'Weekend', 'is_day', 'is_morning'  
, 'is_night', 'same_str']  
  
/opt/conda/lib/python3.6/site-packages/lightgbm/basic.py:762: UserWarning:  
categorical_feature in param dict is overridden.  
  
Score Validation : 21.657467051090165  
0bl | 0 / 15 [0:07:08<] 217.8s, best loss: ?]  
  
/opt/conda/lib/python3.6/site-packages/lightgbm/basic.py:1209: UserWarning:  
categorical_feature in Dataset is overridden.  
New categorical feature is ['City_Atlanta', 'City_Boston', 'City_Chicago', 'City_Philadelphia', 'Ent  
ryType', 'ExitType', 'Hour', 'Intersec', 'IntersectionId', 'Month', 'Weekend', 'is_day', 'is_morning'  
, 'is_night', 'same_str']  
  
/opt/conda/lib/python3.6/site-packages/lightgbm/basic.py:762: UserWarning:  
categorical_feature in param dict is overridden.  
  
Score Validation : 25.91830442123546  
0bl | 0 / 15 [0:09:18<] 217.8s, best loss: ?]  
  
/opt/conda/lib/python3.6/site-packages/lightgbm/basic.py:1209: UserWarning:  
categorical_feature in Dataset is overridden.  
New categorical feature is ['City_Atlanta', 'City_Boston', 'City_Chicago', 'City_Philadelphia', 'Ent  
ryType', 'ExitType', 'Hour', 'Intersec', 'IntersectionId', 'Month', 'Weekend', 'is_day', 'is_morning'  
, 'is_night', 'same_str']  
  
/opt/conda/lib/python3.6/site-packages/lightgbm/basic.py:762: UserWarning:  
categorical_feature in param dict is overridden.  
  
Score Validation : 65.14003527847642  
0bl | 0 / 15 [1:11:30<] 717.8s, best loss: ?]  
  
/opt/conda/lib/python3.6/site-packages/lightgbm/basic.py:1209: UserWarning:  
categorical_feature in Dataset is overridden.  
New categorical feature is ['City_Atlanta', 'City_Boston', 'City_Chicago', 'City_Philadelphia', 'Ent  
ryType', 'ExitType', 'Hour', 'Intersec', 'IntersectionId', 'Month', 'Weekend', 'is_day', 'is_morning'  
, 'is_night', 'same_str']  
  
/opt/conda/lib/python3.6/site-packages/lightgbm/basic.py:762: UserWarning:  
categorical_feature in param dict is overridden.  
  
Score Validation : 6.564807428395332  
7Bl | 1 / 15 [1:15:29<] 3:11:24, 820.30s/it, best loss: 69.31777186872799]  
  
/opt/conda/lib/python3.6/site-packages/lightgbm/basic.py:1209: UserWarning:  
categorical_feature in Dataset is overridden.  
New categorical feature is ['City_Atlanta', 'City_Boston', 'City_Chicago', 'City_Philadelphia', 'Ent  
ryType', 'ExitType', 'Hour', 'Intersec', 'IntersectionId', 'Month', 'Weekend', 'is_day', 'is_morning'  
, 'is_night', 'same_str']  
  
/opt/conda/lib/python3.6/site-packages/lightgbm/basic.py:762: UserWarning:  
categorical_feature in param dict is overridden.  
  
Score Validation : 12.97189877852345  
7Bl | 1 / 15 [1:17:03<] 3:11:24, 820.30s/it, best loss: 69.31777186872799]  
  
/opt/conda/lib/python3.6/site-packages/lightgbm/basic.py:1209: UserWarning:  
categorical_feature in Dataset is overridden.  
New categorical feature is ['City_Atlanta', 'City_Boston', 'City_Chicago', 'City_Philadelphia', 'Ent  
ryType', 'ExitType', 'Hour', 'Intersec', 'IntersectionId', 'Month', 'Weekend', 'is_day', 'is_morning'  
, 'is_night', 'same_str']  
  
/opt/conda/lib/python3.6/site-packages/lightgbm/basic.py:762: UserWarning:  
categorical_feature in param dict is overridden.  
  
Score Validation : 22.007111496703793  
7Bl | 1 / 15 [1:15:38<] 3:11:24, 820.30s/it, best loss: 69.31777186872799]  
  
/opt/conda/lib/python3.6/site-packages/lightgbm/basic.py:1209: UserWarning:  
categorical_feature in Dataset is overridden.  
New categorical feature is ['City_Atlanta', 'City_Boston', 'City_Chicago', 'City_Philadelphia', 'Ent  
ryType', 'ExitType', 'Hour', 'Intersec', 'IntersectionId', 'Month', 'Weekend', 'is_day', 'is_morning'  
, 'is_night', 'same_str']  
  
/opt/conda/lib/python3.6/site-packages/lightgbm/basic.py:762: UserWarning:  
categorical_feature in param dict is overridden.  
  
Score Validation : 25.753810892649778  
7Bl | 1 / 15 [2:11:38<] 3:11:24, 820.30s/it, best loss: 69.31777186872799]  
  
/opt/conda/lib/python3.6/site-packages/lightgbm/basic.py:1209: UserWarning:  
categorical_feature in Dataset is overridden.  
New categorical feature is ['City_Atlanta', 'City_Boston', 'City_Chicago', 'City_Philadelphia', 'Ent  
ryType', 'ExitType', 'Hour', 'Intersec', 'IntersectionId', 'Month', 'Weekend', 'is_day', 'is_morning'  
, 'is_night', 'same_str']  
  
/opt/conda/lib/python3.6/site-packages/lightgbm/basic.py:762: UserWarning:  
categorical_feature in param dict is overridden.  
  
Score Validation : 75.05141893767532  
7Bl | 1 / 15 [2:15:59<] 3:11:24, 820.30s/it, best loss: 69.31777186872799]  
  
/opt/conda/lib/python3.6/site-packages/lightgbm/basic.py:1209: UserWarning:  
categorical_feature in Dataset is overridden.  
New categorical feature is ['City_Atlanta', 'City_Boston', 'City_Chicago', 'City_Philadelphia', 'Ent  
ryType', 'ExitType', 'Hour', 'Intersec', 'IntersectionId', 'Month', 'Weekend', 'is_day', 'is_morning'  
, 'is_night', 'same_str']  
  
/opt/conda/lib/python3.6/site-packages/lightgbm/basic.py:762: UserWarning:  
categorical_feature in param dict is overridden.  
  
Full Score Run : 151.63390301314774  
Parameters: {'boosting': 'gbdt', 'coarsample_bytree': 0.7883586397777255, 'gpu_device_id': 0, 'learning_ra  
te': 0.045508240433874, 'max_depth': 12, 'metric': 'rmse', 'mini_child_samples': 100, 'num_leaves':  
60, 'objective': 'regression', 'subsample': 0.7}  
7Bl | 2 / 15 [2:45:48<] 4:47:49, 774.58s/it, best loss: 69.31777186872799]
```







```
[53]: %time
import lightgbm as lgb

for i in range(len(all_preds)):
    print(f'## {i+1} Run')
    X_tr,X_val,y_tr,y_val=train_test_split(X, y.iloc[:,1],
                                          test_size=0.10, random_state=31)

    xg_train = lgb.Dataset(X_tr, label = y_tr)
    xg_valid = lgb.Dataset(X_val, label = y_val )

    lgbm_reg = lgb.train(best_params, xg_train, 10000,
                        valid_sets = [xg_valid],
                        verbose_eval=500,
                        early_stopping_rounds = 250)

    all_preds[i] = lgbm_reg.predict(X_test, num_iteration=lgbm_reg.best_iteration)

print(f'##{i+1} running done.* ')

## 1 Run
Training until validation scores don't improve for 250 rounds.
[500] valid_0's rmse: 5.9966
[1000] valid_0's rmse: 5.91271
[1500] valid_0's rmse: 5.71554
[2000] valid_0's rmse: 5.66125
[2500] valid_0's rmse: 5.62021
[3000] valid_0's rmse: 5.59185
[3500] valid_0's rmse: 5.5684
[4000] valid_0's rmse: 5.55066
[4500] valid_0's rmse: 5.5384
[5000] valid_0's rmse: 5.52517
[5500] valid_0's rmse: 5.51543
[6000] valid_0's rmse: 5.50762
[6500] valid_0's rmse: 5.50118
[7000] valid_0's rmse: 5.49581
[7500] valid_0's rmse: 5.49113
[8000] valid_0's rmse: 5.48716
[8500] valid_0's rmse: 5.48506
[9000] valid_0's rmse: 5.48223
[9500] valid_0's rmse: 5.47914
[10000] valid_0's rmse: 5.47726
Did not meet early stopping. Best iteration is:
(9578) valid_0's rmse: 5.47704
1 running done.

## 2 Run
Training until validation scores don't improve for 250 rounds.
[500] valid_0's rmse: 11.7668
[1000] valid_0's rmse: 11.0558
[1500] valid_0's rmse: 10.6517
[2000] valid_0's rmse: 10.4188
[2500] valid_0's rmse: 10.2593
[3000] valid_0's rmse: 10.1452
[3500] valid_0's rmse: 10.0319
[4000] valid_0's rmse: 9.97555
[4500] valid_0's rmse: 9.91836
[5000] valid_0's rmse: 9.86238
[5500] valid_0's rmse: 9.81686
[6000] valid_0's rmse: 9.77614
[6500] valid_0's rmse: 9.74299
[7000] valid_0's rmse: 9.71049
[7500] valid_0's rmse: 9.68545
[8000] valid_0's rmse: 9.658
[8500] valid_0's rmse: 9.6307
[9000] valid_0's rmse: 9.62057
[9500] valid_0's rmse: 9.60493
[10000] valid_0's rmse: 9.59146
Did not meet early stopping. Best iteration is:
(10000) valid_0's rmse: 9.59146
2 running done.

## 3 Run
Training until validation scores don't improve for 250 rounds.
[500] valid_0's rmse: 20.4292
[1000] valid_0's rmse: 19.0196
[1500] valid_0's rmse: 18.2787
[2000] valid_0's rmse: 17.8637
[2500] valid_0's rmse: 17.5385
[3000] valid_0's rmse: 17.305
[3500] valid_0's rmse: 17.1265
[4000] valid_0's rmse: 16.9832
[4500] valid_0's rmse: 16.859
[5000] valid_0's rmse: 16.757
[5500] valid_0's rmse: 16.6741
[6000] valid_0's rmse: 16.5996
[6500] valid_0's rmse: 16.5386
[7000] valid_0's rmse: 16.477
[7500] valid_0's rmse: 16.4278
[8000] valid_0's rmse: 16.3848
[8500] valid_0's rmse: 16.3496
[9000] valid_0's rmse: 16.3175
[9500] valid_0's rmse: 16.2871
[10000] valid_0's rmse: 16.2591
Did not meet early stopping. Best iteration is:
(10000) valid_0's rmse: 16.2591
3 running done.

## 4 Run
Training until validation scores don't improve for 250 rounds.
[500] valid_0's rmse: 23.0849
[1000] valid_0's rmse: 22.5089
[1500] valid_0's rmse: 22.253
[2000] valid_0's rmse: 22.0976
[2500] valid_0's rmse: 21.9956
[3000] valid_0's rmse: 21.924
[3500] valid_0's rmse: 21.8876
[4000] valid_0's rmse: 21.845
[4500] valid_0's rmse: 21.82
[5000] valid_0's rmse: 21.8023
[5500] valid_0's rmse: 21.79
[6000] valid_0's rmse: 21.7766
[6500] valid_0's rmse: 21.7734
Early stopping, best iteration is:
(6671) valid_0's rmse: 21.7674
4 running done.

## 5 Run
Training until validation scores don't improve for 250 rounds.
[500] valid_0's rmse: 53.3206
[1000] valid_0's rmse: 50.5741
[1500] valid_0's rmse: 49.3754
[2000] valid_0's rmse: 48.5289
[2500] valid_0's rmse: 48.1257
[3000] valid_0's rmse: 47.7898
[3500] valid_0's rmse: 47.5662
[4000] valid_0's rmse: 47.3932
[4500] valid_0's rmse: 47.2312
[5000] valid_0's rmse: 47.1279
[5500] valid_0's rmse: 47.0027
[6000] valid_0's rmse: 46.9357
[6500] valid_0's rmse: 46.887
[7000] valid_0's rmse: 46.8448
[7500] valid_0's rmse: 46.8243
[8000] valid_0's rmse: 46.7942
[8500] valid_0's rmse: 46.7559
Early stopping, best iteration is:
(8570) valid_0's rmse: 46.7514
5 running done.

## 6 Run
Training until validation scores don't improve for 250 rounds.
[500] valid_0's rmse: 97.7012
[1000] valid_0's rmse: 90.1463
[1500] valid_0's rmse: 87.1521
[2000] valid_0's rmse: 85.2443
[2500] valid_0's rmse: 84.1762
[3000] valid_0's rmse: 83.3879
[3500] valid_0's rmse: 82.9284
[4000] valid_0's rmse: 82.5356
[4500] valid_0's rmse: 82.2854
[5000] valid_0's rmse: 82.0698
[5500] valid_0's rmse: 81.8858
[6000] valid_0's rmse: 81.745
[6500] valid_0's rmse: 81.6304
[7000] valid_0's rmse: 81.533
[7500] valid_0's rmse: 81.4556
[8000] valid_0's rmse: 81.4098
[8500] valid_0's rmse: 81.3495
[9000] valid_0's rmse: 81.3104
[9500] valid_0's rmse: 81.2524
[10000] valid_0's rmse: 81.2331
Did not meet early stopping. Best iteration is:
(9831) valid_0's rmse: 81.2285
6 running done.

CPU times: user 6h 28min 53s, sys: 2min 15s, total: 6h 31min 8s
Wall time: 3h 22min 42s
```

### Importing submission file

- stacking all results in the same file

```
In [54]: sub = pd.read_csv("../input/bigquery-geotab-intersection-congestion/sample_submission.csv")

In [55]: dt = pd.DataFrame(all_preds).stack()
dt = pd.DataFrame(dt)
sub['Target'] = dt['0'].values

In [56]: sub.head()

Out[56]:
   Target0  Target
0      0.0  0.402081
1      0.1  4.877809
2      0.2  11.635321
3      0.3  -5.298662
4      0.4  35.841463

In [57]: sub.to_csv("lgbm_pred_hyperopt_test.csv", index = False)
```

Most part of the first modeling try I got from @danofar  
Please visit the kernel with all work here: <https://www.kaggle.com/danofar/baseline-feature-engineering-geotab-69-5-1b>  
The Catboost model I got from @rohitpatil kernel Link: <https://www.kaggle.com/rohitpatil/geotab-catboost>  
Some ideas of modelling I saw on: <https://www.kaggle.com/dcaichara/feature-engineering-and-lightgbm>

**NOTE: This Kernel is not finished.**

**Please stay tuned and votes up the kernel, please!**