

# CLUSTERING GBO LOCATION DATA INTELLIGENTLY IN PYTHON

We have taxi rank locations. Our task is to determine where we can build our service centers/<sup>stations</sup> to accommodate the taxis in the region.

Simple way  $\Rightarrow$  Make clusters of data and build a service center/<sup>station</sup> at centroid of each cluster.

Import Libraries and <sup>set</sup> colours

**TASK 1: Exploratory Data Analysis** and <sup>Removing NULLS</sup> Duplicates.

```
df = pd.read_csv('Data/taxi-data.csv')
```

```
df.head()
```

Column name of duplicate value we need to check

Duplicate

```
df.duplicated(subset=['LONG', 'LAT']).values.any()
```

<sup>return true if any is true.</sup>

Is none or empty

```
df.isna().values.any()
```

```
df.shape()  $\rightarrow$  Before dropping data.
```

```
df.dropna(inplace=True)
```

$\rightarrow$  Make changes to same data instead of returning new data

```
df.drop_duplicates(subset=['LONG', 'LAT'], keep='first', inplace=True)
```



df.shape() → after dropping

x = np.array(df[['LAT', 'LON']], dtype='float64')

x → array only of LON & LAT.

plt.scatter(x[:,0], x[:,1], alpha=0.2, s=50)

everything from first column      everything from second column      size of circle notes everything light for better view in dense conditions.

## TASK -2 VISUALISE GEOGRAPHICAL DATA

We need to do it so we can zoom in data to properly find clusters. Point where our map should start.

m = folium.Map(location=[df.LAT.mean(), df.LON.mean()], zoom\_start=9, tiles='Stamen Jones', colour='scheme')

m → It will print only the map. we need to populate it. the way we do is iteratively add all rows in data frame.

for \_, row in df.iterrows():  
we use underscore because iterrows sends row number and then data of row. we do not need the number and hence we ignore it.



Making circle for all rows.

folium.CircleMarker(

location = [row.LAT, row.LON],

radius = 5,

popup = re.sub(r'[^a-zA-Z ]+', '', row.Name)

colour = '#1787FE'

fill\_colour = 'blue',

fill\_colour = '#1787FE'

) .add\_to(m)

add to map

m → Now the map contains circle for each taxi

Because we have not actually cleaned the strings in row name it is possible there might be some form of escapes that cause problem when you convert the python code to javascript or html file.

Way to Deal → Take row name and only keep characters and spaces. One way to do that is to use regex.

re.sub(r'[^a-zA-Z ]+', '', row.Name)

↑  
substitute reg  
exp  
chars

Anything except  
these will be  
replaced by



## TASK-3 Clustering Strength / Performance Metric

# Get dummy data to work with

$X$ -blobs,  $\rightarrow$  make-blobs ( $n$ -samples = 1000,  
array  $\rightarrow$  to ignore  $\rightarrow$  total samples

10 Clusters  $\rightarrow$  centers = 10,  
 $x$  &  $y$  values  $\rightarrow$   $n$ -features = 2,  
standard deviation  $\rightarrow$  cluster-std = 0.5,  
random-state = 4)

plt.scatter( $X$ -blobs[:, 0],  $X$ -blobs[:, 1], alpha=0.2)  
Algo to cluster

class-predictions = np.load('Data / Sample-clusters.npy')

class-predictions  $\rightarrow$  value of each cluster number

unique-clusters = np.unique(class-predictions)  
for unique-cluster in unique-clusters:

$X = X$ -blobs[class-predictions == unique-cluster]  
plt.scatter( $X$ [:, 0],  $X$ [:, 1], alpha=0.2,  
 $c = \text{cols}[\text{unique-cluster}]$ )

$\rightarrow$  It gives score from -1 to +1, better  
 $\downarrow$  Score better clusters.

array of colors  
before task 1.

Silhouette-score ( $X$ -blobs, class-predictions)  
Data  $\rightarrow$  Cluster Labels

If we use sample-clusters-improved.npy  
we get better score.



K-Means algo starts with random values. and can result in different.

## TASK-4 K-MEANS CLUSTERING

$x = \text{np.array}(df[['lon', 'LAT']], \text{dtype} = \text{'float64'})$   
 $k=70$  # Any number not defined here we can have no. of service station we can actually build.

fitting to model.

$\text{model} = \text{KMeans}(n\text{-clusters} = k, \text{random\_state} = 11)$   
 $\text{fit}(x)$

predicting clusters

$\text{class\_predictions} = \text{model.predict}(x)$

making new column in our data

$df['cluster\_k\text{-means}\{k\}'] = \text{class\_predictions}$

$df.head()$  - viewing the new column.

new dataframe (column we care about)

$\text{def create\_map}(df, \text{cluster\_column})$ :

starting location

$m = \text{folium.Map}(\text{location} = [df.LAT.mean(), df.LON.mean()])$

starting zoom

$\text{zoom\_start} = 9$ ,

our colour files = 'Stamen Toner'  
scheme.

for row in df.iter\_rows():

$\text{cluster\_colour} = \text{cols}[\text{row}[\text{cluster\_column}]]$

T

Colours according to cluster label.



forum, CircleMaker(  
locations = [row['LAT'], row['Lon']],  
radius = 5,

What will show → popup = row[cluster - column],  
in popup = the clusters  
label. colour = cluster - column colour,  
fill = True,  
fill\_color = cluster - colour

) .add = to (m).

defined  
on peer page.

return m.

The new column  
name.

m = create\_map(df, CLUSTER\_kmeans, 70)

m = print map

Calculate best silhouette score.

best = best - silhouette, best\_k = -1, 0

Initial values → value -1 0

for k in tqdm(range(2, 100)):

2-100 clusters → model = K-means(n\_clusters = k

random\_state = 1)  
.fit(x)

class\_predictions = model.predict(x)

Current score → cur\_silhouette = silhouette\_score(x, class\_predictions)

If cur\_silhouette > best\_silhouette,

best\_k = k

best\_silhouette = cur\_silhouette



```
print (best_k)
print (best_silhouette)
```

## TASK-5 DBSCAN

Density Based Spatial Clustering of Applications with Noise.

K-means doesn't take in account density. Like more clusters in dense regions and less on outside but here we know in city centers there will be dense clusters and hence more clusters

dummy with -1 as noise.

```
dummy = np.array([-1, -1, -1, 2, 3, 4, 5, -1])
```

make each noise  
in individual clusters

```
new = np.array([counter + 2] * x if x == -1 else x for counter, x in enumerate(dummy))
```

```
dummy = [-2, -3, -4, 2, 3, 4, 5, -6]
```

radius as min. each point min samples in clusters.

```
model = DBSCAN(eps = 0.01, min_samples = 5).fit(x)
```

class predictions = model.labels - # predict won't work.

```
df['Clusters - DBSCAN'] = class_predictions
```

```
m = create_map(df, 'CLUSTERS - DBSCAN')
```

## TASK-6 HDBSCAN <sup>Heirarchical</sup> → Hybrid DBSCAN

DBSCAN have global understanding of density in dataset. This means it can't necessarily pickup different kinds of Density within a specific dataset.



HDBSCAN will try to learn multiple types of densities.

DBSCAN takes out outliers. That is not a downside but given our dataset problem we want each point to be assigned a cluster.

```
model = hdbscan.HDBSCAN (min-cluster-size=5,  
                          min-samples=2,  
                          cluster-selection-epsilon=0.01)
```

```
class_predictions = model.fit_predict(X)  
df['CLUSTER_HDBSCAN'] = class_predictions
```

```
m = create_map(df, 'CLUSTER_HDBSCAN')
```

**TASK-7 ADDRESSING OUTLIERS**  $K=NN \rightarrow K$  nearest  
classifier = KNeighborsClassifier (n\_neighbors=1) neighbors

Training Set  
Prediction Set  
 $\rightarrow$  df\_train = df[df.CLUSTER\_HDBSCAN != -1] Identified Clusters  
 $\rightarrow$  df\_predict = df[df.CLUSTER\_HDBSCAN == -1] Outliers

```
X_train = np.array(df_train[['LONG', 'LAT']], dtype='float64')  
y_train = np.array(df_train['CLUSTER_HDBSCAN'], dtype='float64')  
X_predict = np.array(df_predict[['LONG', 'LAT']], dtype='float64')
```

```
classifier = KNeighborsClassifier()  
classifier.fit(X_train, y_train)  $\rightarrow$  fit the model  
predictions = classifier.predict(X_predict)  $\rightarrow$  predict values
```

Copy values to new column  
replace  
df['CLUSTER\_hybrid'] = df['CLUSTER\_HDBSCAN']

```
df.loc[df.CLUSTER_HDBSCAN == -1, 'CLUSTER_hybrid'] =  
      predictions
```

where it is -1      Replace clusters hybrid column with predictions

```
m = create_map(df, 'CLUSTER_hybrid')
```