

Clustering Geolocation Data Intelligently in Python

We have taxi rank locations, and want to define key clusters of these taxis where we can build service stations for all taxis operating in that region.

Project Outline

[Task 1](#): Exploratory Data Analysis

[Task 2](#): Visualizing Geographical Data

[Task 3](#): Clustering Strength / Performance Metric

[Task 4](#): K-Means Clustering

[Task 5](#): DBSCAN

```
In [3]: 1 import matplotlib
2 %matplotlib inline
3 %config InlineBackend.figure_format = 'svg'
4 import matplotlib.pyplot as plt
5 plt.style.use('ggplot')
6
7 import pandas as pd
8 import numpy as np
9
10 from tqdm import tqdm
11
12 from sklearn.cluster import KMeans, DBSCAN
13 from sklearn.metrics import silhouette_score
14 from sklearn.datasets import make_blobs
15 from sklearn.neighbors import KNeighborsClassifier
16
17 from ipywidgets import interactive
18
19 from collections import defaultdict
20
21 #import hdbscan
22 import folium
23 import re
24
25
26 cols = ['#e6194b', '#3cb44b', '#ffe119', '#4363d8', '#f58231', '#911eb4',
27         '#46f0f0', '#f032e6', '#bcf60c', '#fabebe', '#008080', '#e6beff',
28         '#9a6324', '#fffac8', '#800000', '#aaffc3', '#808000', '#ffd8b1',
29         '#000075', '#808080']*10
```

Task 1: Exploratory Data Analysis

```
In [4]: 1 df = pd.read_csv('Data/taxi_data.csv')
```

```
In [5]: 1 df.head()
```

Out[5]:

	LON	LAT	NAME
0	28.17858	-25.73882	11th Street Taxi Rank
1	28.17660	-25.73795	81 Bazaar Street Taxi Rank
2	27.83239	-26.53722	Adams Road Taxi Rank
3	28.12514	-26.26666	Alberton City Mall Taxi Rank
4	28.10144	-26.10567	Alexandra Main Taxi Rank

```
In [6]: 1 df.duplicated(subset=['LON', 'LAT']).values.any()
```

Out[6]: True

```
In [7]: 1 df.isna().values.any()
```

Out[7]: True

```
In [8]: 1 print(f'Before dropping NaNs and dupes\t:\tdf.shape = {df.shape}')
2 df.dropna(inplace=True)
3 df.drop_duplicates(subset=['LON', 'LAT'], keep='first', inplace=True)
4 print(f'After dropping NaNs and dupes\t:\tdf.shape = {df.shape}')
```

```
Before dropping NaNs and dupes : df.shape = (838, 3)
After dropping NaNs and dupes : df.shape = (823, 3)
```

```
In [9]: 1 df.head()
```

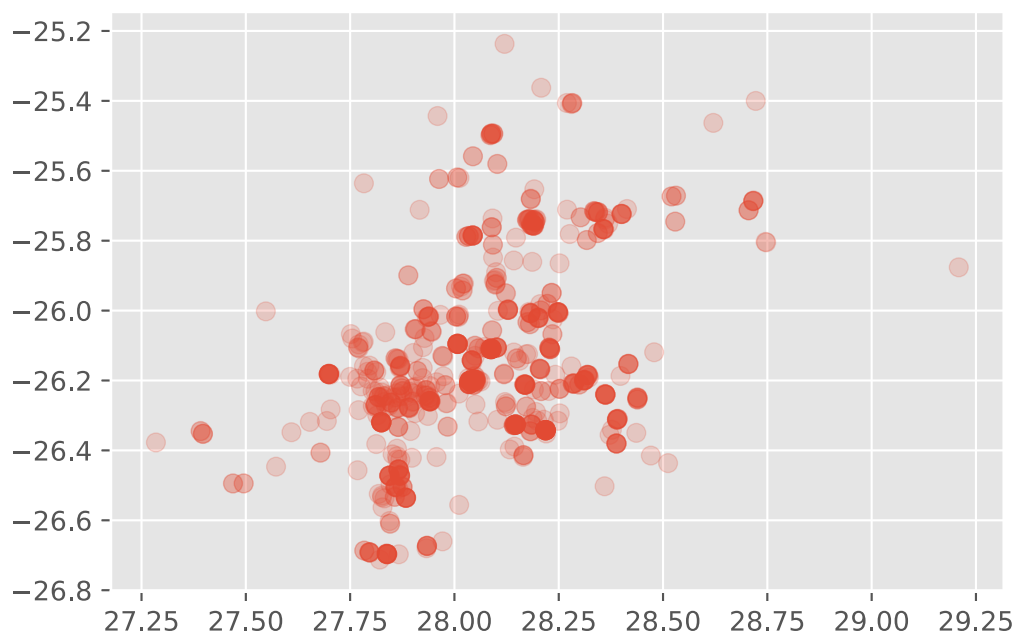
Out[9]:

	LON	LAT	NAME
0	28.17858	-25.73882	11th Street Taxi Rank
1	28.17660	-25.73795	81 Bazaar Street Taxi Rank
2	27.83239	-26.53722	Adams Road Taxi Rank
3	28.12514	-26.26666	Alberton City Mall Taxi Rank
4	28.10144	-26.10567	Alexandra Main Taxi Rank

```
In [10]: 1 X = np.array(df[['LON', 'LAT']], dtype='float64')
```

```
In [11]: 1 plt.scatter(X[:,0], X[:,1], alpha=0.2, s=50)
```

```
Out[11]: <matplotlib.collections.PathCollection at 0x21d2090b6a0>
```

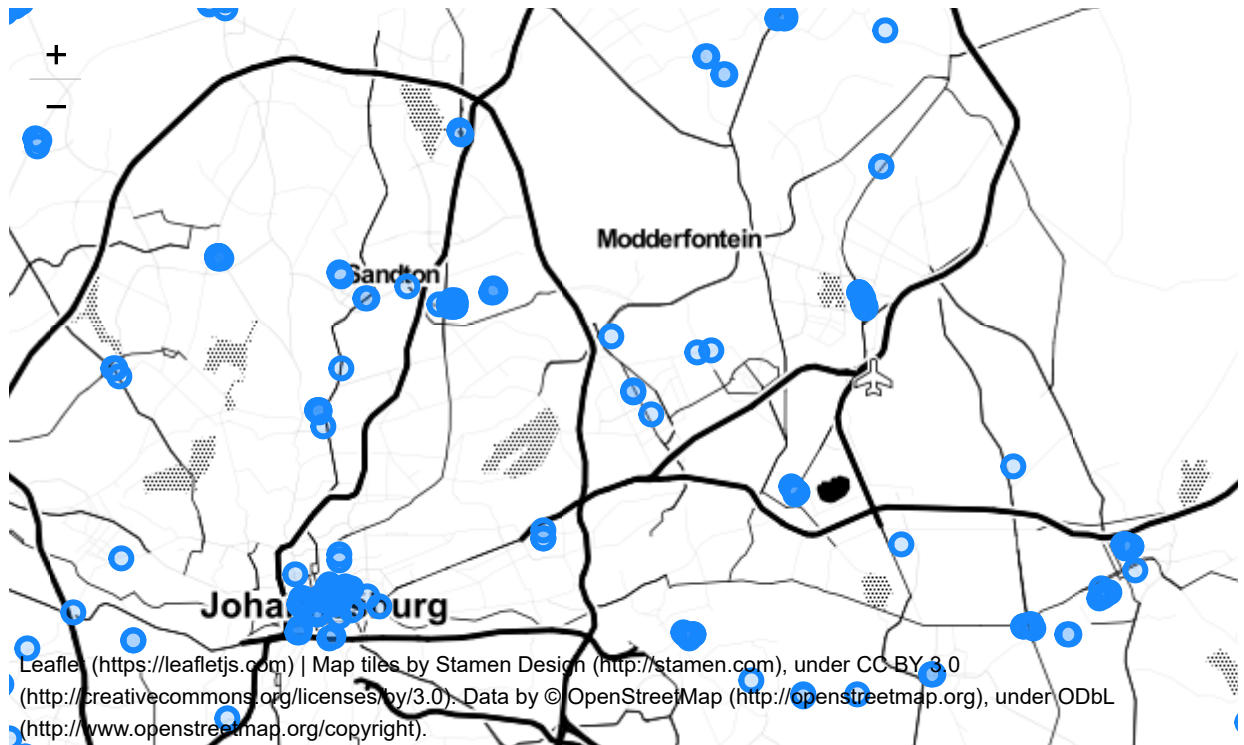


Task 2: Visualizing Geographical Data

```
In [12]: 1 m = folium.Map(location=[df.LAT.mean(), df.LON.mean()], zoom_start=9,
2           tiles='Stamen Toner')
3
4 for _, row in df.iterrows():
5     folium.CircleMarker(
6         location=[row.LAT, row.LON],
7         radius=5,
8         popup=re.sub(r'^a-zA-Z ]+', '', row.NAME),
9         color='#1787FE',
10        fill=True,
11        fill_colour='#1787FE'
12    ).add_to(m)
```

In [13]: 1 m

Out[13]:

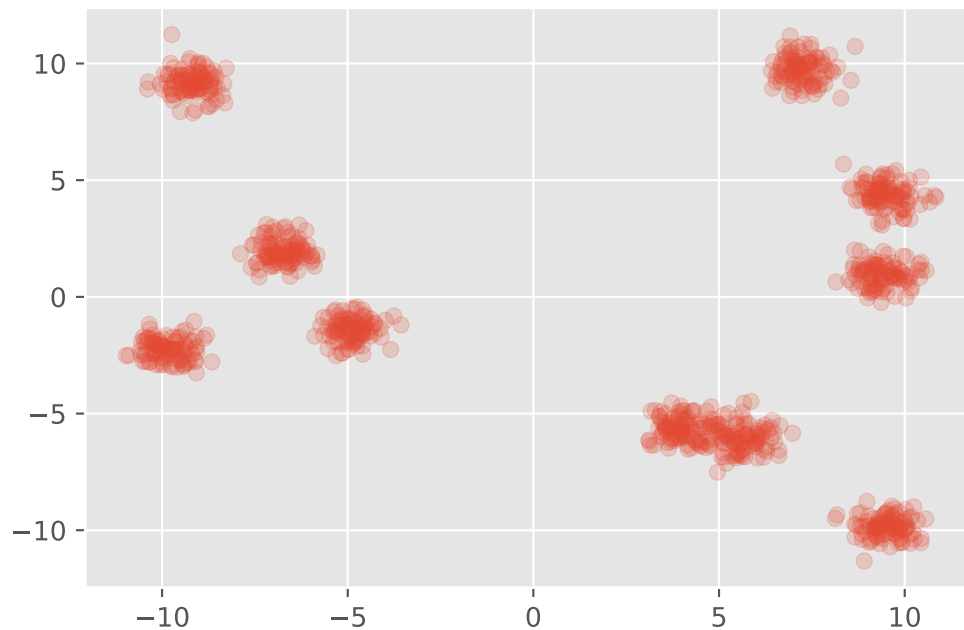


Task 3: Clustering Strength / Performance Metric

In [14]: 1 X_blobs, _ = make_blobs(n_samples=1000, centers=10, n_features=2,
2 cluster_std=0.5, random_state=4)

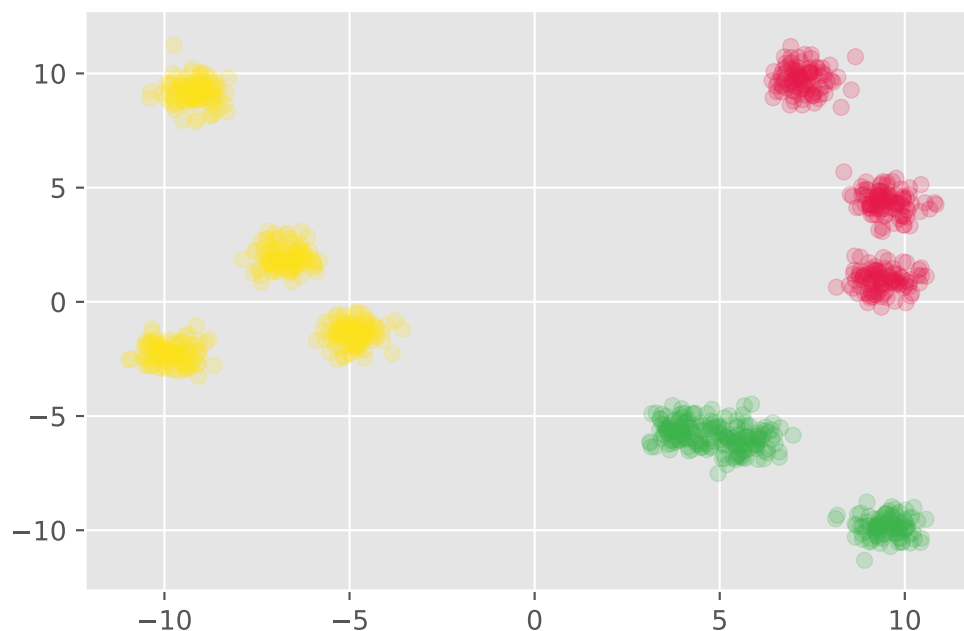
```
In [15]: 1 plt.scatter(X_blobs[:,0], X_blobs[:,1], alpha=0.2)
```

```
Out[15]: <matplotlib.collections.PathCollection at 0x21d2147f320>
```



```
In [16]: 1 class_predictions = np.load('Data/sample_clusters.npy')
```

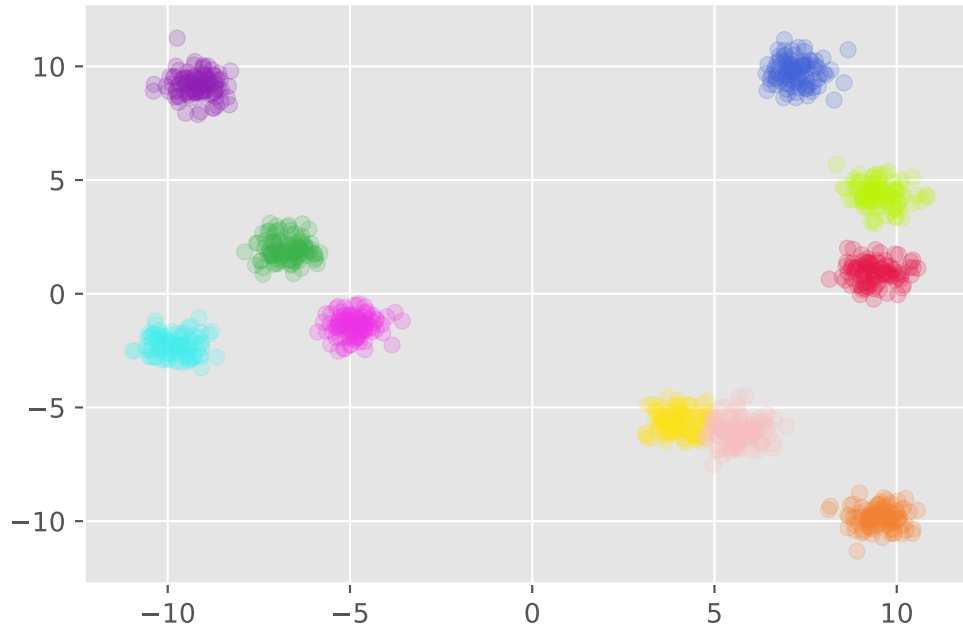
```
In [17]: 1 unique_clusters = np.unique(class_predictions)
2 for unique_cluster in unique_clusters:
3     X = X_blobs[class_predictions==unique_cluster]
4     plt.scatter(X[:,0], X[:,1], alpha=0.2, c=cols[unique_cluster])
```



```
In [18]: 1 silhouette_score(X_blobs, class_predictions)
```

```
Out[18]: 0.6657220862867241
```

```
In [19]: 1 class_predictions = np.load('Data/sample_clusters_improved.npy')
2 unique_clusters = np.unique(class_predictions)
3 for unique_cluster in unique_clusters:
4     X = X_blobs[class_predictions==unique_cluster]
5     plt.scatter(X[:,0], X[:,1], alpha=0.2, c=cols[unique_cluster])
```



```
In [20]: 1 silhouette_score(X_blobs, class_predictions)
```

```
Out[20]: 0.7473587799908298
```

Task 4: K-Means Clustering

```
In [21]: 1 X_blobs, _ = make_blobs(n_samples=1000, centers=50,
2                               n_features=2, cluster_std=1, random_state=4)
```

```
In [22]: 1 data = defaultdict(dict)
2 for x in range(1,21):
3     model = KMeans(n_clusters=3, random_state=17,
4                   max_iter=x, n_init=1).fit(X_blobs)
5
6     data[x]['class_predictions'] = model.predict(X_blobs)
7     data[x]['centroids'] = model.cluster_centers_
8     data[x]['unique_classes'] = np.unique(class_predictions)
```

```

In [23]: 1 def f(x):
2         class_predictions = data[x]['class_predictions']
3         centroids = data[x]['centroids']
4         unique_classes = data[x]['unique_classes']
5
6         for unique_class in unique_classes:
7             plt.scatter(X_blobs[class_predictions==unique_class][:,0],
8                         X_blobs[class_predictions==unique_class][:,1],
9                         alpha=0.3, c=cols[unique_class])
10        plt.scatter(centroids[:,0], centroids[:,1], s=200, c='#000000', marker='x')
11        plt.ylim([-15,15]); plt.xlim([-15,15])
12        plt.title('How K-Means Clusters')
13
14        interactive_plot = interactive(f, x=(1, 20))
15        output = interactive_plot.children[-1]
16        output.layout.height = '350px'
17        interactive_plot

```

A Jupyter widget could not be displayed because the widget state could not be found. This could happen if the kernel storing the widget is no longer available, or if the widget state was not saved in the notebook. You may be able to create the widget by running the appropriate cells.

```

In [24]: 1 X = np.array(df[['LON', 'LAT']], dtype='float64')
2         k = 70
3         model = KMeans(n_clusters=k, random_state=17).fit(X)
4         class_predictions = model.predict(X)
5         df[f'CLUSTER_kmeans{k}'] = class_predictions

```

```

In [25]: 1 df.head()

```

Out[25]:

	LON	LAT	NAME	CLUSTER_kmeans70
0	28.17858	-25.73882	11th Street Taxi Rank	1
1	28.17660	-25.73795	81 Bazaar Street Taxi Rank	1
2	27.83239	-26.53722	Adams Road Taxi Rank	9
3	28.12514	-26.26666	Alberton City Mall Taxi Rank	8
4	28.10144	-26.10567	Alexandra Main Taxi Rank	4

```

In [26]: 1 def create_map(df, cluster_column):
2         m = folium.Map(location=[df.LAT.mean(), df.LON.mean()], zoom_start=9, ti
3
4         for _, row in df.iterrows():
5
6             if row[cluster_column] == -1:
7                 cluster_colour = '#000000'
8             else:
9                 cluster_colour = cols[row[cluster_column]]
10
11             folium.CircleMarker(
12                 location= [row['LAT'], row['LON']],
13                 radius=5,
14                 popup= row[cluster_column],
15                 color=cluster_colour,
16                 fill=True,
17                 fill_color=cluster_colour
18             ).add_to(m)
19
20         return m
21
22 m = create_map(df, 'CLUSTER_kmeans70')
23 print(f'K={k}')
24 print(f'Silhouette Score: {silhouette_score(X, class_predictions)}')
25
26 m.save('kmeans_70.html')

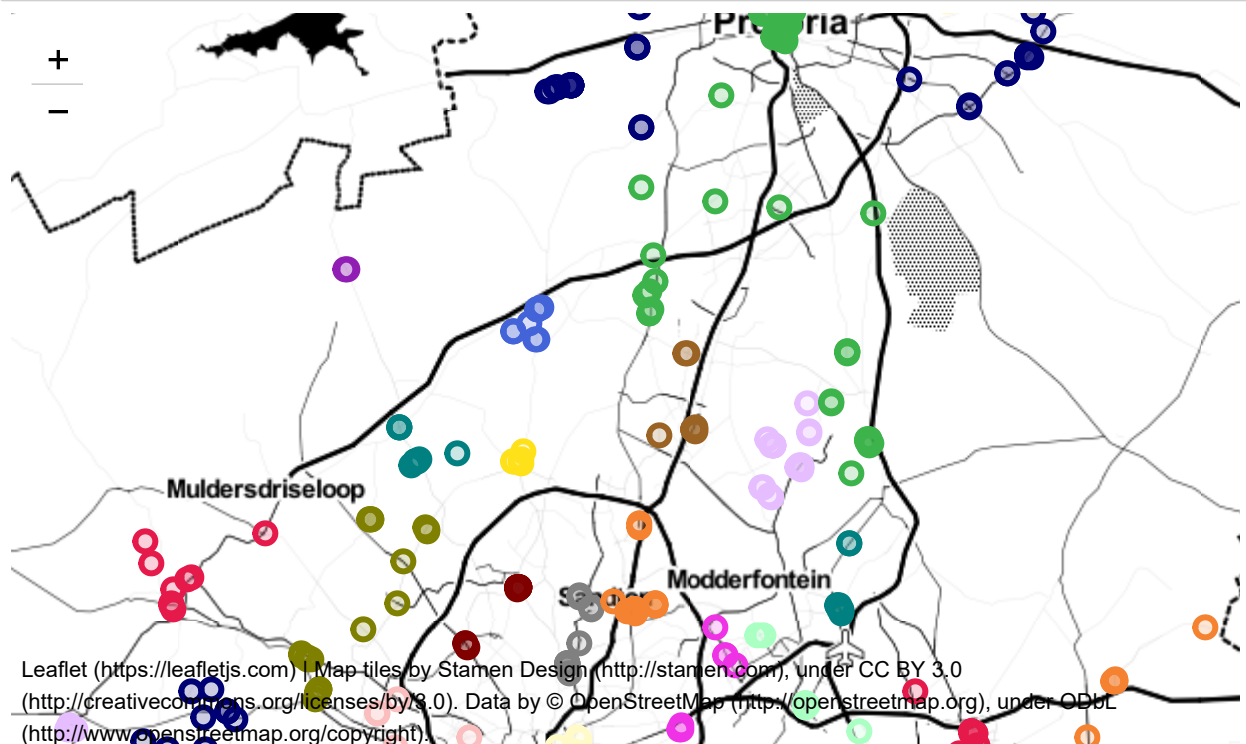
```

K=70

Silhouette Score: 0.6527069281188838

In [27]: 1 m

Out[27]:




```
In [28]: 1 best_silhouette, best_k = -1, 0
2
3 for k in tqdm(range(2, 100)):
4     model = KMeans(n_clusters=k, random_state=1).fit(X)
5     class_predictions = model.predict(X)
6
7     curr_silhouette = silhouette_score(X, class_predictions)
8     if curr_silhouette > best_silhouette:
9         best_k = k
10        best_silhouette = curr_silhouette
11
12 print(f'K={best_k}')
13 print(f'Silhouette Score: {best_silhouette}')
```

```

In [31]: 1 m = create_map(df, 'CLUSTERS_DBSCAN')
          2
          3
          4 print(f'Number of clusters found: {len(np.unique(class_predictions))}')
          5 print(f'Number of outliers found: {len(class_predictions[class_predictions==
          6
          7 print(f'Silhouette ignoring outliers: {silhouette_score(X[class_predictions!
          8
          9 no_outliers = 0
         10 no_outliers = np.array([(counter+2)*x if x==-1 else x for counter, x in enum
         11 print(f'Silhouette outliers as singletons: {silhouette_score(X, no_outliers)

```

Number of clusters found: 51
 Number of outliers found: 289
 Silhouette ignoring outliers: 0.9232138250288208
 Silhouette outliers as singletons: 0.5667489350583482

```

In [32]: 1 m

```

Out[32]:

