

DBSCAN

In this notebook you will use GPU-accelerated DBSCAN to identify clusters of infected people.

Objectives

By the time you complete this notebook you will be able to:

- Use GPU-accelerated DBSCAN
- Use cuXfilter to visualize DBSCAN clusters

Imports

```
In [1]: import cudf
import cuml

import cuxfilter as cxf
```

Load Data

For this notebook, we again load a subset of our population data with only the columns we need. An `infected` column has been added to the data to indicate whether or not a person is known to be infected with our simulated virus.

```
In [2]: gdf = cudf.read_csv('./data/pop_2-04.csv', dtype=['float32', 'float32', 'float32'])
print(gdf.dtypes)
gdf.shape
```

```
northing    float32
easting     float32
infected     float32
dtype: object
(1000000, 3)
```

Out[2]:

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

```
Out[3]:
```

	northing	easting	infected
0	178547.296875	368012.1250	0.0
1	174068.281250	543802.1250	0.0
2	358293.687500	435639.8750	0.0
3	87240.304688	389607.3750	0.0
4	158261.015625	340764.9375	0.0

```
In [4]: gdf['infected'].value_counts()
```

```
Out[4]: 0.0    984331
        1.0    15669
        Name: infected, dtype: int32
```

DBSCAN Clustering

DBSCAN is another unsupervised clustering algorithm that is particularly effective when the number of clusters is not known up front and the clusters may have concave or other unusual shapes--a situation that often applies in geospatial analytics.

In this series of exercises you will use DBSCAN to identify clusters of infected people by location, which may help us identify groups becoming infected from common patient zeroes and assist in response planning.

Exercise: Make a DBSCAN Instance

Create a DBSCAN instance by using `cuml.DBSCAN`. Pass in the named argument `eps` (the maximum distance a point can be from the nearest point in a cluster to be considered possibly in that cluster) to be `5000`. Since the `northing` and `easting` values we created are measured in meters, this will allow us to identify clusters of infected people where individuals may be separated from the rest of the cluster by up to 5 kilometers.

```
In [8]: dbscan = cuml.DBSCAN(eps=5000)
```

Solution

```
In [9]: # %Load solutions/dbscan_instance
        dbscan = cuml.DBSCAN(eps=5000)
```

Exercise: Identify Infected Clusters

Create a new dataframe from rows of the original dataframe where `infected` is `1` (true), and call it `infected_df` --be sure to reset the dataframe's index afterward. Use `dbscan.fit_predict` to perform clustering on the `northing` and `easting` columns of `infected_df`, and turn the resulting series into a new column in `infected_gdf` called "cluster". Finally, compute the number of clusters identified by DBSCAN.

```
In [11]: infected_df = gdf[gdf['infected'] == 1].reset_index()
        infected_df['cluster'] = dbscan.fit_predict(infected_df[['northing', 'easting']])
        infected_df['cluster'].nunique()
```

```
Out[11]: 96
```

Solution

```
In [12]: # %Load solutions/identify_infected
infected_df = gdf[gdf['infected'] == 1].reset_index()
infected_df['cluster'] = dbscan.fit_predict(infected_df[['northing', 'easting']])
infected_df['cluster'].nunique()
```

Out[12]: 96

Visualize the Clusters

Because we have the same column names as in the K-means example-- `easting`, `northing`, and `cluster`--we can use the same code to visualize the clusters.

Associate a Data Source with cuXfilter

```
In [13]: cxf_data = cxf.DataFrame.from_dataframe(infected_df)
```

Define Charts and Widgets

As in the K-means notebook, we have an existing integer column to use with multi-select:
`cluster`.

```
In [14]: chart_width = 600
scatter_chart = cxf.charts.datashader.scatter(x='easting', y='northing',
                                              width=chart_width,
                                              height=int((gdf['northing'].max() - gdf[
                                                (gdf['easting'].max() - gdf[
                                                  chart_width))

cluster_widget = cxf.charts.panel_widgets.multi_select('cluster')
```

Create and Show the Dashboard

```
In [15]: dash = cxf_data.dashboard(charts=[scatter_chart], sidebar=[cluster_widget], theme=cxf.t
```

```
In [16]: scatter_chart.view()
```

Out[16]:

```
In [17]: %%js
var host = window.location.host;
element.innerText = "'http://" + host + "'";
```

Set `my_url` in the next cell to the value just printed, making sure to include the quotes:

```
In [18]: my_url = 'http://dli-604a4aa51b37-32d463.aws.labs.courses.nvidia.com'  
dash.show(my_url, port=8789)
```

Dashboard running at port 8789

Out[18]:

... and you can run the next cell to generate a link to the dashboard:

```
In [19]: %%js  
var host = window.location.host;  
var url = 'http://' + host + '/lab/proxy/8789/';  
element.innerHTML = '<a style="color:blue;" target="_blank" href='+url+'>Open Dashboard
```

```
In [20]: dash.stop()
```

Please Restart the Kernel

```
In [ ]: import IPython  
app = IPython.Application.instance()  
app.kernel.do_shutdown(True)
```

Next

In the next notebook, you will use GPU-accelerated logistic regression to estimate infection risk based on features of our population members.