

Multi-GPU K-Means with Dask

In this notebook you will use GPU-accelerated K-means to identify population clusters in a multi-node, multi-GPU scalable way with Dask.

Objectives

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

- Use distributed, GPU-accelerated K-means with Dask

Imports

First we import the needed modules to create a Dask cuDF cluster.

```
In [ ]: import subprocess

from dask.distributed import Client, wait, progress
from dask_cuda import LocalCUDACluster
```

After that, we create the cluster.

```
In [ ]: cmd = "hostname --all-ip-addresses"
process = subprocess.Popen(cmd.split(), stdout=subprocess.PIPE)
output, error = process.communicate()
IPADDR = str(output.decode()).split()[0]

cluster = LocalCUDACluster(ip=IPADDR)
client = Client(cluster)
```

Finally, as we did before, we import CUDA context creators after setting up the cluster so they don't lock to a single device.

```
In [ ]: import cudf
import dask_cudf

import cuml
from cuml.dask.cluster import KMeans
```

Load and Persist Data

We will begin by loading the data, The dataset has the two grid coordinate columns, `easting` and `northing`, derived from the main population dataset we have prepared.

```
In [ ]: ddf = dask_cudf.read_csv('./data/pop5x_2-07.csv', dtype=['float32', 'float32'])
```

Training the K-means model is very similar to both the scikit-learn version and the cuML single-GPU version--by setting up the client and importing from the `cuml.dask.cluster` module, the algorithm will automatically use the local Dask cluster we have set up.

Note that calling `.fit` triggers Dask computation.

```
In [ ]: %%time
        dkm = KMeans(n_clusters=20)
        dkm.fit(ddf)
```

Once we have the fit model, we extract the cluster centers and rename the columns from their generic '0' and '1' to reflect the data on which they were trained.

```
In [ ]: cluster_centers = dkm.cluster_centers_
        cluster_centers.columns = ddf.columns
        cluster_centers.dtypes
```

Exercise: Count Members of the Southernmost Cluster

Using the `cluster_centers`, identify which cluster is the southernmost (has the lowest `northing` value) with the `nsmallest` method, then use `dkm.predict` to get labels for the data, and finally filter the labels to determine how many individuals the model estimated were in that cluster.

```
In [1]: %load solutions/southernmost_cluster
```

Please Restart the Kernel

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

Next

In [the next notebook](#), you will calculate infection risk again, this time using the powerful XGBoost algorithm.