

## DDA\_exercise3.py

```

import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from mpi4py import MPI

# initialize MPI
comm = MPI.COMM_WORLD
rank = comm.Get_rank()
size = comm.Get_size()

#load data form csv file into pandas dataframe called pandas
#specify column header is the first row in data file
# index_col 0 specifies that the first column in the CSV file should be used as the
DataFrame index.
pandas = pd.read_csv("cluster_data.csv", header=0, index_col=0)

plt.scatter(pandas['x'], pandas['y'])
plt.xlabel('X')
plt.ylabel('Y')
plt.show()

# split data into equal chinks
chunk_size = len(pandas) // size
chunk_start = rank * chunk_size
chunk_end = chunk_start + chunk_size
if rank == size - 1:
    chunk_end = len(pandas)
chunk = pandas[chunk_start:chunk_end]

# initialize the centroids
centroids = np.random.uniform(low = 0, high=10, size=(4, 2))

# Assign points to the closest centroids
distance = np.sqrt(((chunk.values - centroids[:, np.newaxis])**2).sum(axis=2))
assignment = np.argmin(distance, axis=0)

# compute the local centroids for this MPI Process
local_centroids = np.zeros((4, 2))
counts = np.zeros(4)
for i in range(len(chunk)):
    local_centroids[assignment[i]] += chunk.iloc[i][['x', 'y']]
    counts[assignment[i]] += 1
local_centroids /= np.where(counts == 0, 1e-9, counts)[:, np.newaxis]

#Reduce the local centroids to the global centroids using mpi
global_centroids = np.zeros((4, 2))
global_counts = np.zeros(4)
for i in range(4):
    comm.Reduce(local_centroids[i], global_centroids[i], op=MPI.SUM, root=0)
    comm.Reduce(counts[i], counts, op=MPI.SUM, root=0)
#Divide the global centroids by the number of assigned points to get the avg centroids
for i in range(4):

```

```
    if global_counts[i] > 0:  
        global_centroids /= global_counts[i]  
  
#print the global centroids from the coordinator(rank 0)  
if rank == 0:  
    print("Global centroids:", global_centroids)
```