# Distributed Data Analytics

## *Exercise Sheet 3*

Mohsan Jameel, Florian Pal

Information Systems and Machine Learning Lab

University of Hildesheim

Name: Pedram Babakhani

Student Number: 276848

Email: Pedram_Babakhani@yahoo.com

# Exercise 1.

Clustering is the task of assigning a set of objects into groups (called clusters) so that the objects in the same cluster are more similar (in some sense or another) to each other than to those in other clusters. As a naïve example: animals can be clustered as land animals, water animals and amphibians.

• k-means clustering is a method of clustering which aims to partition n data points into k clusters (n >> k) in which each observation belongs to the cluster with the nearest mean.

• The nearness is calculated by distance function which is mostly Euclidian distance or Manhattan distance.(I useed Euclidian)

 • One important assumption to be made is the data points are independent of each other. In other words there exists no dependency between any data points.

In the following you can see the time execution comparison between parallel and sequential approach.
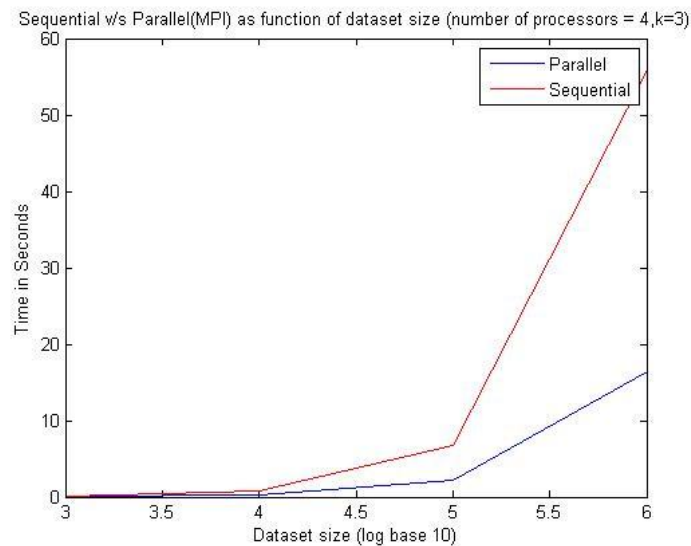


Figure 1: Time execution between parallel and sequential[1]

In the following, figure 2 demonstrate the speedup and time execution vs number of processes, we can see as we increase number of processes we will get better speedup but we have to find the best K and best number of Processor.
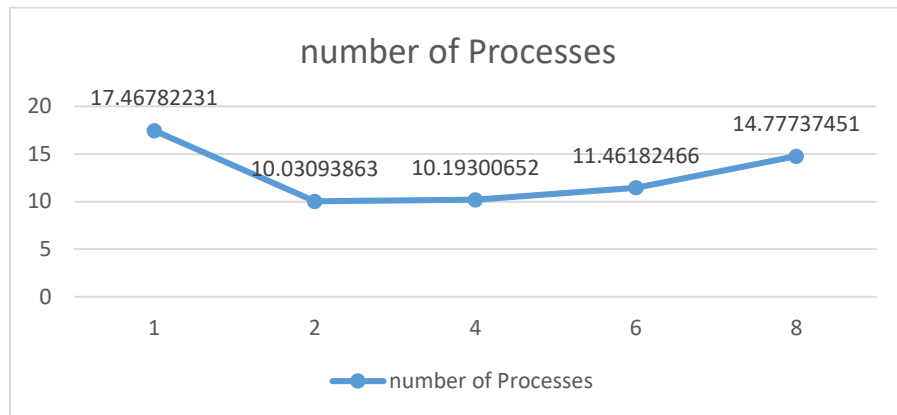
---

[1] Buffalo university
https://www.cse.buffalo.edu/faculty/miller/Courses/CSE633/Chandramohan-Fall-2012-CSE633.pdf

Figure 2: Time execution and speed up vs number of processes[2]

First of all I victories the texts in folders (only 2 folders because of memory limitation), After that I split the vector horizontally and send for all processes, every process try to find Euclidian Distance and then find which of them are closer the update the membership and assign the data to a cluster. I used dictionary for membership, actually it is dictionary of list which keeps the membership of every class. On the other hand at the end I will calculate new centroids by mean of the values in membership vector.

Then at the end, process 0 receives the membership from other processes and calculate new centroids and calculate time execution. In the following section I will compare my results with the real results.

Figure 3 shows the output of terminal.



Figure 3: Terminal output

[2] Buffalo university
https://www.cse.buffalo.edu/faculty/miller/Courses/CSE633/Chandramohan-Fall-2012-CSE633.pdf

3

## Exercise 2.

I calculate time by process zero, because it takes longer than others because it has to do split the data then apply k-means and gather data from others. Graph 1 shows the execution time simulated by k=4 and different number of processes.

As we can see 2 processes will be the best number of processes when k=4 in k means algorithm.



number of Processes

| | |
|---|---|
| 17.46782231 | |
| 10.03093863 | 10.19300652 |
| 11.46182466 | 14.77737451 |

Graph 1: execution time vs number of processor

Generally, when we increase number of processor it does not mean that we are decreasing time execution, sometimes synchronization between processes takes longer than the task we are operating parallel then my result sounds to be true because in figure 2 (left part) we can see that always more processor does not decrease the time execution.

Note: all time execution is simulated by my code, you can run my code and check the result. Furtheremore, I have put comment enough in my code Hope it would be enough.