

# APPM4058A & APPM7045A Project: K-Means Clustering & Fuzzy C-Means Clustering

1846346, 1669326, COMS Hons

June 27, 2021

WITS  
UNIVERSITY



## Problem Domain

Fuzzy C-Means Clustering and K-Means clustering are algorithms that are typically used in semi-supervised & unsupervised Machine Learning. K-Means clustering aims to segment  $N$  observations into  $K$  clusters where each observation belongs to the cluster with the nearest mean. It aims to minimise within cluster variance. K-Means clustering (Dunn [1973]) can be used for cluster analysis or feature learning. The problem is considered NP-Hard. Fuzzy C-Means clustering (Macqueen [1967]) is considered a soft version of K-Means clustering where the observation has associations to multiple clusters rather than just one. Fuzzy C-Means clustering can be used in image analysis and marketing. Given that these algorithms are computationally expensive & computation can be conducted in parallel. This makes the use of parallel programming techniques highly desirable. Thus, we will explore the effects of using CUDA<sup>1</sup> & Open MPI<sup>2</sup> frameworks on computation time.

Definitions The following definitions are relevant in this report:

- Observation: A data point in  $d$ -dimensional space that has an associated cluster/s.
- K-Mean: A reference point that is designated to be a  $d$ -dimensional cluster. To which observations are assigned. They represent the average location of their assigned clusters.
- C-Mean: A reference point with some degree of fuzziness. It is a  $d$ -dimensional cluster. All observations will have some level of association  $p$ , where  $p \in [0, 1]$  to the reference point. They represent the average location based on observations' associations.

## Methodology

For each algorithm, we will give a brief overview of the pseudocode. The pseudocode will be implemented in our serial versions. We will highlight the components of the serial code that will be modified to make use of parallel methods for an improved performance.

<sup>1</sup><https://developer.nvidia.com/cuda-gpus>

<sup>2</sup><https://www.open-mpi.org/>

---

### K-Means Clustering

---

- 1: Given  $N$  observations &  $K$  Means
  - 2: Initialise  $K$  Means at random points in  $d$ -dimensional space
  - 3: For  $i=0,1,...,iterations$
  - 4:   Calculate distance to each Mean for each observation
  - 5:   Assign observation to closest Mean
  - 6:   Update the Means
  - 7: Return observations
- 

### Fuzzy C-Means Clustering

---

- 1: Given  $N$  observations &  $C$  Means
  - 2: Initialise  $C$  Means at random points in  $d$ -dimensional space
  - 3: For  $i=0,1,...,iterations$
  - 4:   Calculate distance to each Mean for each observation
  - 5:   Update association of each data point to an observation
  - 6:   Update the C-Means centres
  - 7: Return observations
- 

For K-Means clustering we will parallelise the assignment of observations to each cluster. In CUDA we will aim to assign the maximum number of threads. Using this, each thread will take a single observation and calculate the distance to each cluster. This will update the observations assignment to a cluster. In Open MPI we will aim to apply a similar approach. Here we will distribute the observations to each process such that each process receives  $\frac{\text{Number of Observations}}{\text{Number of Processes}}$ , we will refer to this number as  $A$ . Each process calculates the distances for their  $A$  observations and assigns them to the relevant cluster. From here the averaging of the clusters is done in 2 steps. The first is that every node will calculate a partial sum of where the node should be, the master node then takes the sum of the partial sums and averages it. The new clusters are then broadcasted to each node and the process repeats.

For Fuzzy C-Means clustering we parallelise the calculations in a similar fashion to that of K-Means clustering. In CUDA we will aim to assign the maximum number of threads to update our observations association to various clusters. This requires a distance calculation similar to our K-Means methods however, here we need to make use of

the power function with the exponent being the level of fuzziness. We have set this level of fuzziness to be 4. Thus, this is vastly more computationally intensive compared to the K-Means method. To update the  $C$  clusters we will aim to spawn as many threads as there are clusters, thus limiting us to 1024 clusters. Each thread will update a given cluster using the observations.

In Open MPI we will spawn 8 processes. Each process will calculate  $A$  observations' cluster association. To update our  $C$  clusters we will compute partial sums of the *point sum* & *prob sum*, based on the Fuzzy C-Means algorithm. Each process will then accumulate their results for these partial sums into a single process - the rank 0 process. Using this cumulative result of partial sums we will calculate C-Means updated value.

### Validation of Results

In order to ensure that our parallel methods produce the same results as our serial methods we will do the following for each algorithm:

- Create a serial implementation.
- Initialise our synthetic data.
- Store two sets of the synthetic data at runtime, one set will be used to perform the parallel computations.
- Maintain the order of parallel computations to be the same as that of the serial computations.
- Collect all parallel results into a single result.
- Compute the serial version using the second set of synthetic data.
- Compare each cluster and verify that they are of the same value, or within an acceptable tolerance level.

### Aims

Our aim is to show that a parallel implementation for both the Fuzzy C-Means algorithm & K-Means algorithm provides significant benefits in computational performance. We will showcase the benefits of using CUDA & Open MPI and how these two methods of parallelisation compare to each other.

```

Job is running on node mscluster[10-10]-----
SLURM: sbatch is running on mscluster@ms.wits.ac.za
SLURM: job ID is 131884
SLURM: submit directory is /home-mscluster/jharrisdevey/Projects_2021/HPC/FuzzyMeans
SLURM: number of nodes allocated is 2
SLURM: number of cores is 16
SLURM: job name is FuzzyClustering
-----
CALCULATION FINISHED
success
=====
Numpoints: 1048576, Numcluster: 2, Dimensions: 2
Time for the serial code: 2.33267 SECONDS
Numpoints: 1048576, Numcluster: 2, Dimensions: 2
Time for the kernel: 9.249257 seconds
Test Successful

=====
Numpoints: 1048576, Numcluster: 2, Dimensions: 2
Total time taken 0.184671
Test Successful

=====

```

Figure 1: Proof of Cluster Usage

### Experimental Setup

We performed our computations with computers that had the following noteworthy hardware: Intel Core i9-10940X CPU (14 cores), NVIDIA RTX3090 GPU (24GB), and 128GB of system RAM. We will aim to utilise the maximum number of threads in CUDA implementations. We will aim to use 8 processors for our Open MPI implementations. Computations were performed using High Performance Computing infrastructure provided by the Mathematical Sciences Support unit at the University of the Witwatersrand.

### Evaluation of Results & Discussion

#### Results

##### K-Means

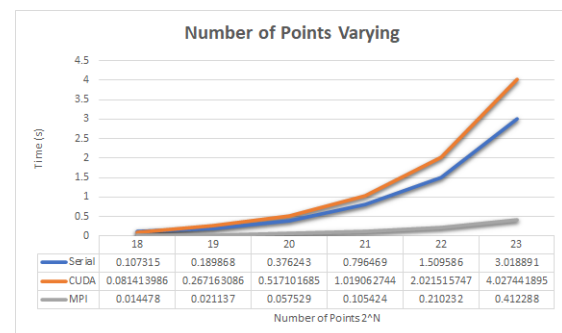
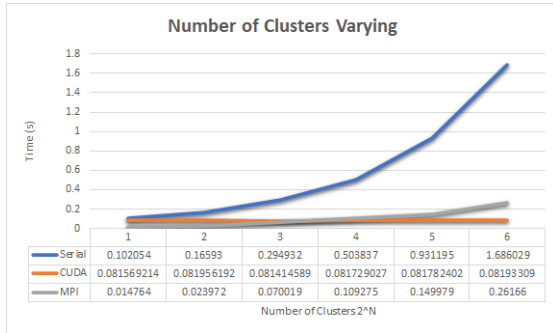
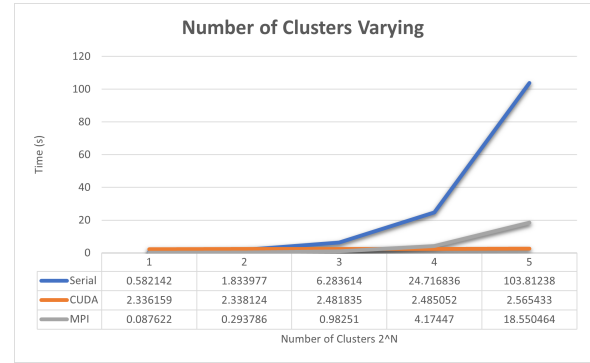
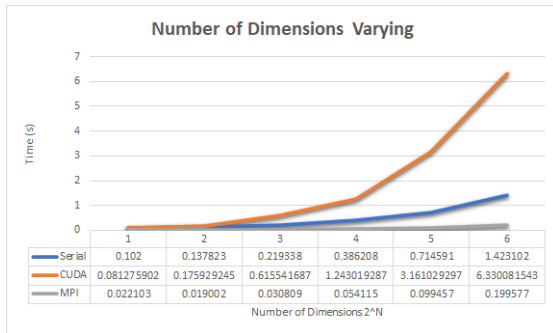
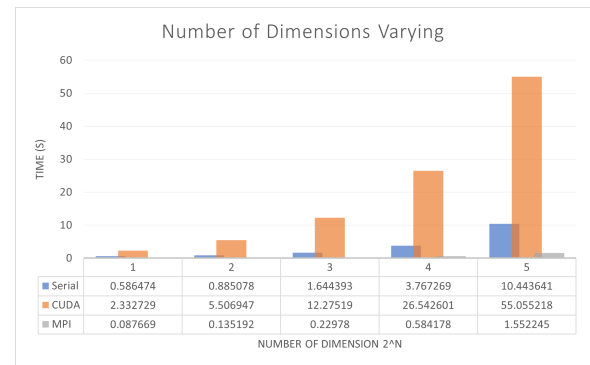
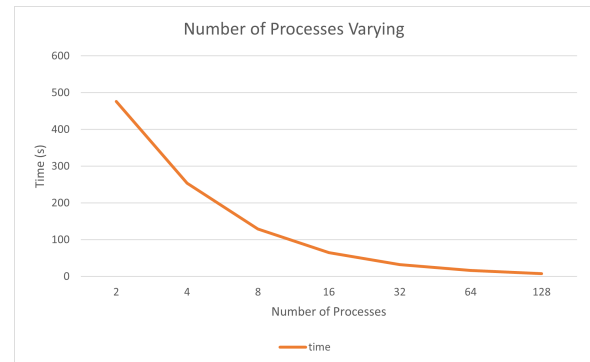
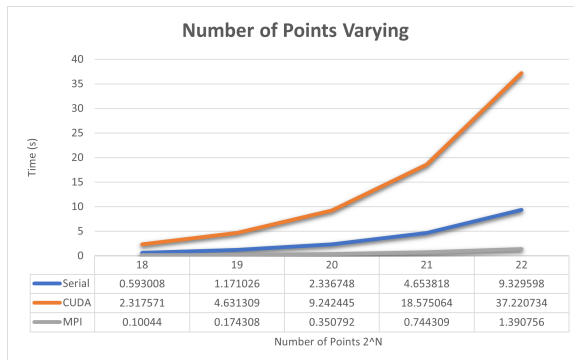


Figure 2: K-Means: Number of Points Varying for 2 Dimensions and 2 Clusters

Figure 3: K-Means: Number of Clusters Varying for 2 Dimensions and  $2^{18}$  PointsFigure 6: Fuzzy C-Means: Number of Clusters Varying for 2 Dimensions and  $2^{18}$  PointsFigure 4: K-Means: Number of Dimensions Varying for 2 Dimensions and  $2^{18}$  PointsFigure 7: Fuzzy C-Means: Number of Dimensions Varying for 2 Dimensions and  $2^{18}$  PointsFigure 8: Fuzzy C-Means: Number of Processes Varying for  $2^3$  Clusters,  $2^3$  Dimensions and  $2^{24}$  PointsFigure 5: Fuzzy C-Means: Number of Clusters Varying for 2 Dimensions and  $2^{18}$  Points

## Fuzzy C-Means

## Discussion of Results

### K-Means

**Varying Points:** Based on the above results (Fig.2) for K-Mean clustering, we can see that as the number of points varies, CUDA gives the worst times, whilst Open MPI gives the best time. The CUDA implementation is based on a

global memory approach and therefore there is a lot of contention as each thread is attempting to access the clusters position. We give a discussion and possible solution about this below in the Limitations & Difficulties section. The Open MPI solution is the best in this case and this makes sense. As the number of points changes the entire problem is being split over multiple nodes. This allows for a massive parallelisation speedup.

**Varying Clusters:** We can see that as we vary the number of clusters Fig.3, the CUDA implementation remains almost constant whilst the serial implementation starts to get exponentially worse. This makes sense since in the CUDA implementation, we assign each cluster to its own thread. This is so that the calculation is happening completely in parallel. This allows us to get a massive parallel advantage compared to the sequential version. The Open MPI version of this also performs well, but is slightly worse than the CUDA version.

**Varying Dimensions:** Finally, for varying the number of dimensions we can see (Fig.4) that CUDA performs the worst here since there is no explicit code to parallelise the loop over the dimensions. So changing the dimensions of the data causes each thread to stay in a for loop for longer. This causes a larger data contention when each thread is accessing the cluster position as it does this more often. Open MPI gives the best results here, as the number of points is split over multiple nodes. The Open MPI code also has no contention over the accessing of a clusters position and hence performs the best.

---

### Fuzzy C-Means

---

**Varying Points:** Based on the above results (Fig.5) for Fuzzy C-Means, we can see that the graphs give a similar output to the K-Means implementation, but just at a longer run time for everything. Fuzzy C-Means is a more expensive calculation than K-Means clustering and hence why the times shown are larger. The overall trend of the graphs remain the same though. The CUDA implementation as shown in the graphs performs the worst, but this is due to the data contention as mentioned before in K-Means. We give a discussion and possible solution about this in the Limitations & Difficulties section below. The Open MPI solution performed the best as the problem was distributed over multiple nodes without the data contention issue.

**Varying Clusters:** We can see that as we vary the number of clusters (Fig.6), the CUDA implementation gives the

best result and remains almost constant. This is because this section of the code was completely parallelisable. We can also see that Open MPI performed well having also beaten the sequential time.

**Varying Dimensions:** Finally, for varying the number of dimensions (Fig.7), the CUDA implementation was the worst due to the data contention whilst Open MPI performed the best. This is a similar scenario to that of K-Means clustering as discussed above.

**Effect of Processors** It can be seen in Fig.8, that by increasing the number of processors we nearly half the computation time. This is highly desirable. This shows that the problem exhibits a high degree of parallelism and benefits greatly from parallel programming methods. Thus using Open MPI with a maximal number of cores for the data set should always be the goal.

### Limitations & Difficulties Encountered

**Floating Points:** We encountered issues with errors compounding fairly quickly when there are a large number of iterations with the random values being in the range of  $10^3$  -  $10^7$ . This makes sense given the extensive use of square roots and power functions in each algorithm. The errors are typically fairly small, often being quite closely related to the chosen maximum value allowed. Thus for very close approximations, our parallel methods work.

**CUDA Memory:** We made extensive use of global memory. This led to a decrease in performance. Given that the calculations for our observations would have at minimum  $2^{18}$  observations, this very quickly leads to bandwidth becoming an issue. This could be improved by storing components of each in constant memory such as; the observation's position & cluster's association. This would greatly reduce the bandwidth used. This could be further optimised to make use of shared memory, which would produce optimal results.

**Processes** Open MPI's performance is directly related to the number of processes that can be spawned. If you have a limited number of processors then Open MPI's results will not generalise. Thus there would be a hard cap on the performance gain that is possible.

## Conclusion

Clearly, it can be seen that parallel methods can be used to dramatically decrease computation time in Fuzzy C-Means & K-Means clustering. Overall the best approach to use is Open MPI. This produced the best results. When there are large numbers of clusters CUDA is also a desirable method. Should the memory be changed for CUDA, this might lead to significantly better performance. Also, a hybrid implementation of the two methods might also generate the best results.

## References

- Dunn, J. (1973). A fuzzy relative of the isodata process and its use in detecting compact well-separated clusters. *Journal of Cybernetics*, 3(3):32–57.
- Macqueen, J. (1967). Some methods for classification and analysis of multivariate observations. *Berkeley Symposium on Mathematical Statistics and Probability*, 5(1):281–297.