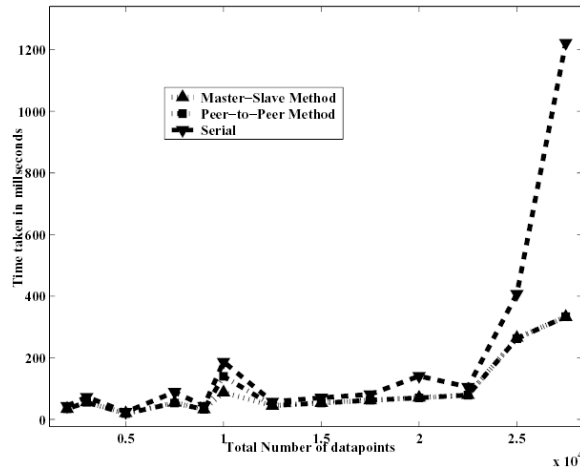


- b. All the processors calculate the centroids for their local data, using Serial KMeans clustering.
- c. All the processors send their local cluster data to rest of the processors.
- d. All the processors receive the data sent by other processors and recompute the centroids locally.
- e. Each processor checks for convergence condition. If convergence condition is not reached, then steps 2 & 3 are repeated. This process is repeated till convergence condition is reached.

Figure 2 compares the two MatlabMPI implementation of K-Means clustering with the Serial implementation. From Fig. 2 it is observed that the difference in the time taken by serial process and that taken by the two MatlabMPI implementations increases as the number of centroids to be clustered or the number of data points to be clustered increases. Moreover, both the parallel implementations take nearly the same amount of time.



This publication was made possible through support provided by DoD HPCMP PET activities through Mississippi State University under the terms of Agreement No. #GS04T01BFC0060. The opinions expressed herein are those of the author(s) and do not necessarily reflect the views of the DoD or Mississippi State University.

References

Vipin Kumar Mahesh V. Joshi, George Karypis [2002]. *Shared memory parallelization of data mining algorithms: Techniques, programming interface, and performanc*. In Second SIAM conference on Data Mining, 2002.

Jeremy Kepner [2002]. MatlabMPI Improves Matlab Performance By 300x. In MAUI HIGH PERFORMANCE COMPUTING CENTER Application Briefs, 2002.