

Performance Comparisons of the K-Means Algorithm

Mark A. Ward

Center for Data Science, NYU

May 19, 2015

Outline

- 1 Introduction
 - The K-Means Algorithm
 - Examples
- 2 Parallel K-Means
 - MPI Communication Overview
 - Initialization
- 3 Performance
- 4 Concluding Remarks

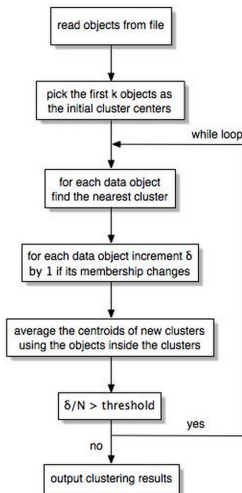
Outline

- 1 Introduction
 - The K-Means Algorithm
 - Examples
- 2 Parallel K-Means
 - MPI Communication Overview
 - Initialization
- 3 Performance
- 4 Concluding Remarks

The K-Means Algorithm

- Algorithm that aims to partition n observations into k clusters where each observation belongs to the cluster nearest mean or centroid
- Vector quantization in signal processing
- Cluster analysis in data mining, how do you choose k ?
- Feature learning in semi-supervised or unsupervised learning
- Clustering is NP-hard, but k-means is a heuristic that converges quickly to a local optimum
- Typical running time $O(nkdi)$, repeat for T trials
- For n points in $[0, 1]^d$ with independent gaussian perturbations with 0 mean and σ^2 variance, upper bound expected running time by $O(n^{34}k^{34}d^8\log^4(n)/\sigma^6)$. [Arthur 2009]

The K-Means Algorithm



N : number of data objects

K : number of clusters

$objects[N]$: array of data objects

$clusters[K]$: array of cluster centers

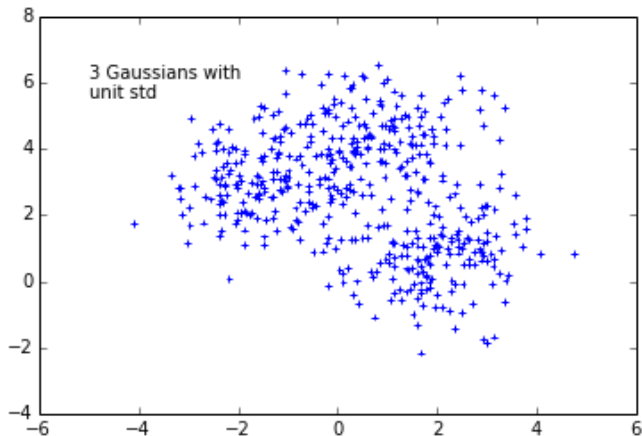
$membership[N]$: array of object memberships

kmeans_clustering()

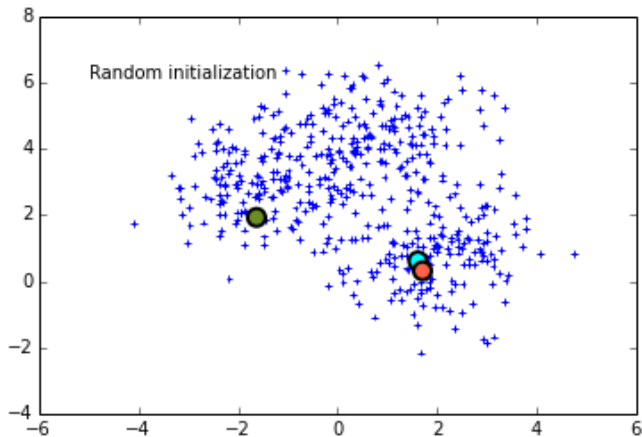
```

1  while  $\delta/N > \text{threshold}$ 
2     $\delta \leftarrow 0$ 
3    for  $i \leftarrow 0$  to  $N-1$ 
4      for  $j \leftarrow 0$  to  $K-1$ 
5         $\text{distance} \leftarrow |objects[i] - clusters[j]|$ 
6        if  $\text{distance} < d_{min}$ 
7           $d_{min} \leftarrow \text{distance}$ 
8           $n \leftarrow j$ 
9        if  $membership[i] \neq n$ 
10          $\delta \leftarrow \delta + 1$ 
11          $membership[i] \leftarrow n$ 
12          $new\_clusters[n] \leftarrow new\_clusters[n] + objects[i]$ 
13          $new\_cluster\_size[n] \leftarrow new\_cluster\_size[n] + 1$ 
14     for  $j \leftarrow 0$  to  $K-1$ 
15        $clusters[j][*] \leftarrow new\_clusters[j][*] / new\_cluster\_size[j]$ 
16        $new\_clusters[j][*] \leftarrow 0$ 
17        $new\_cluster\_size[j] \leftarrow 0$ 
  
```

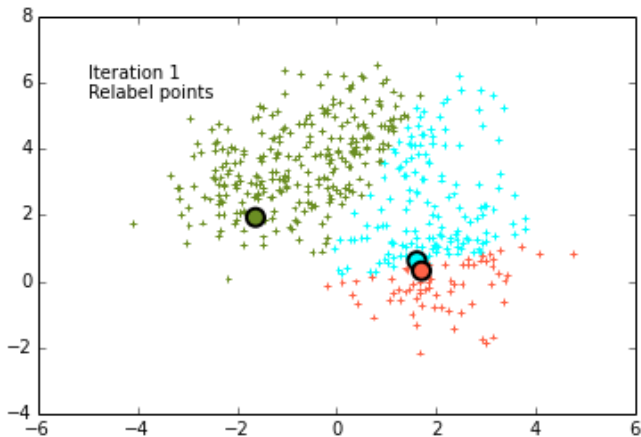
Examples



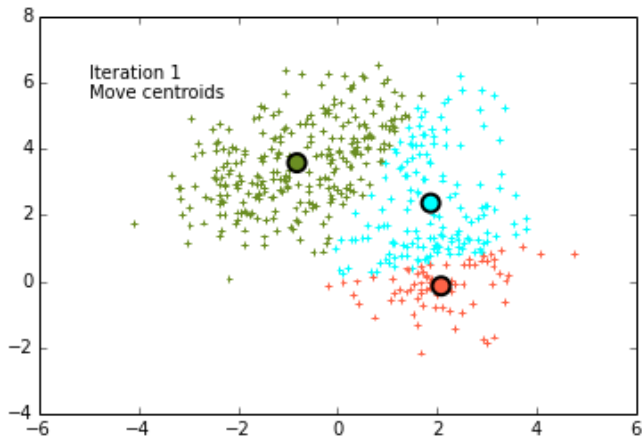
Examples



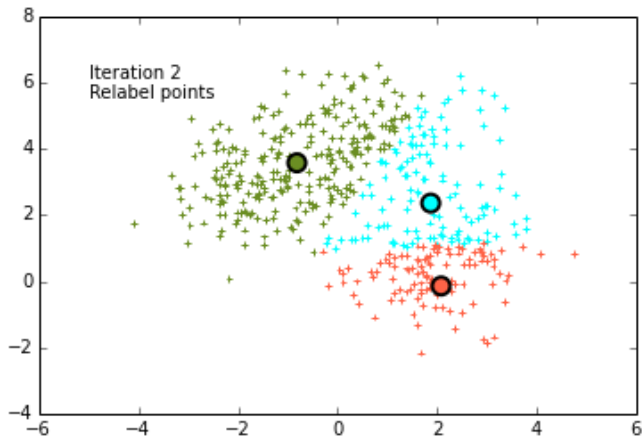
Examples



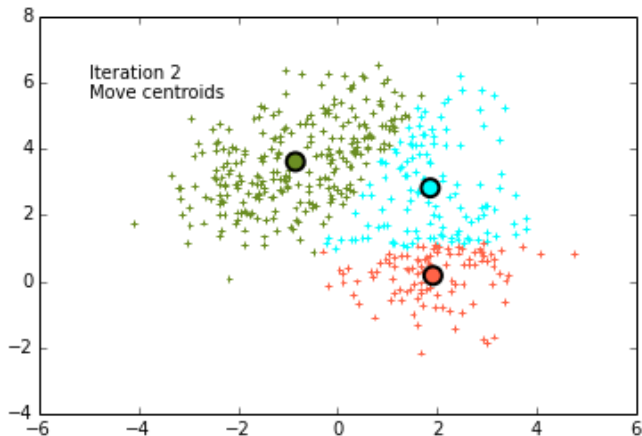
Examples



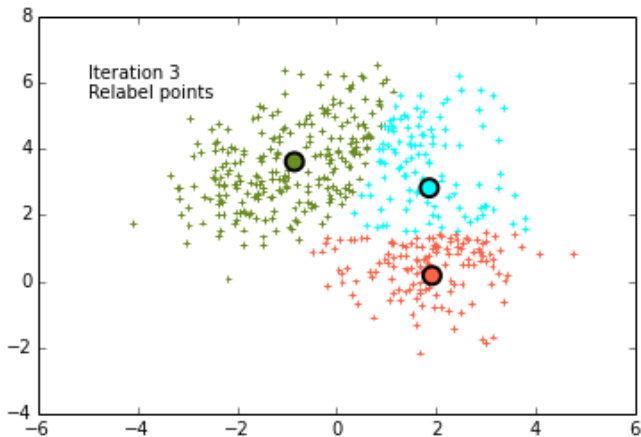
Examples



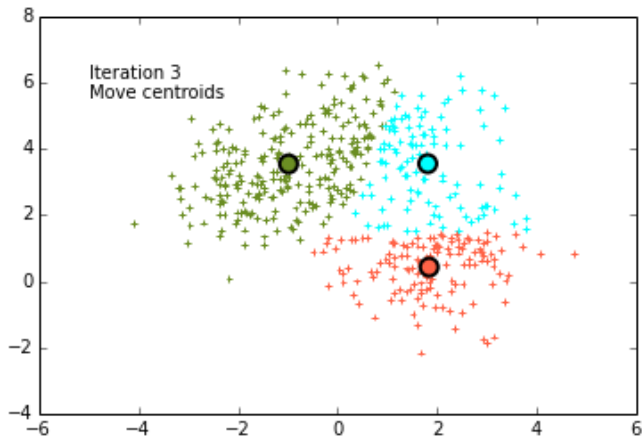
Examples



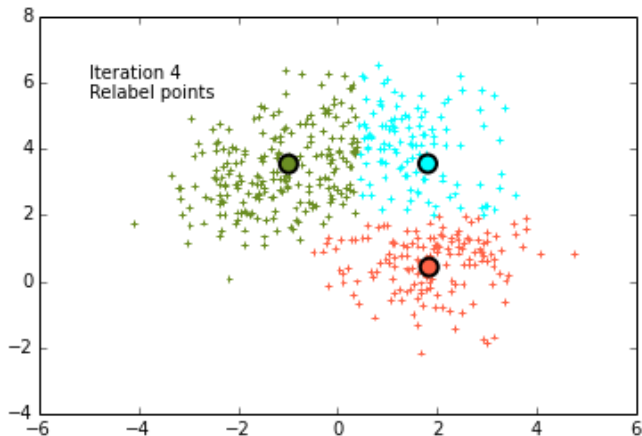
Examples



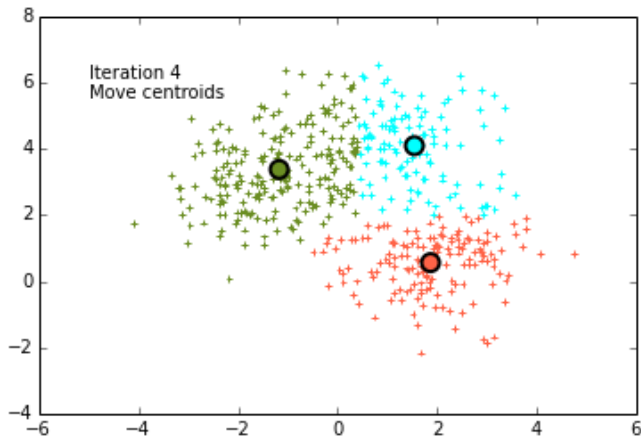
Examples



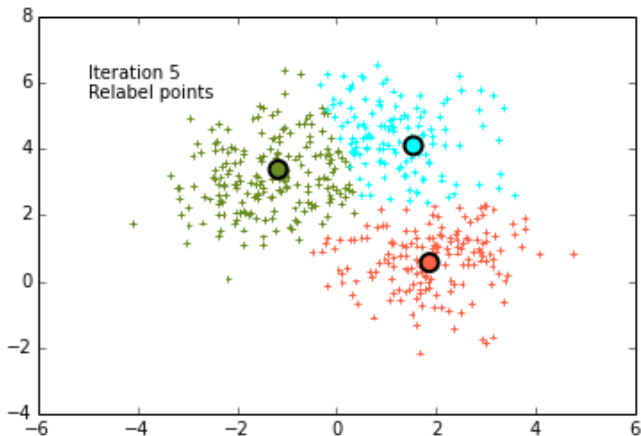
Examples



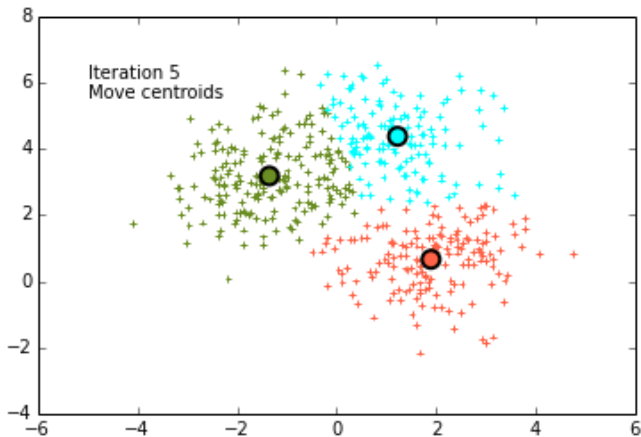
Examples



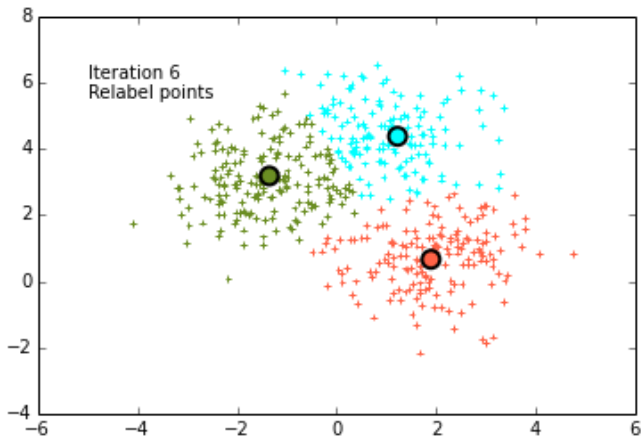
Examples



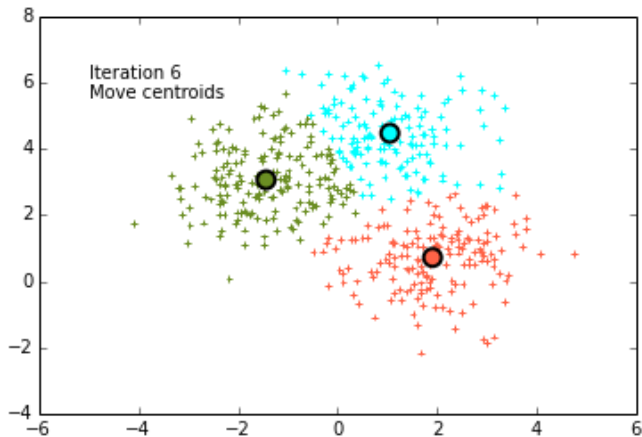
Examples



Examples



Examples



Examples

229,931 colors



Examples

256 colors



Examples

128 colors



Examples

64 colors



Examples

32 colors



Examples

16 colors



Examples

10 colors



Examples

5 colors



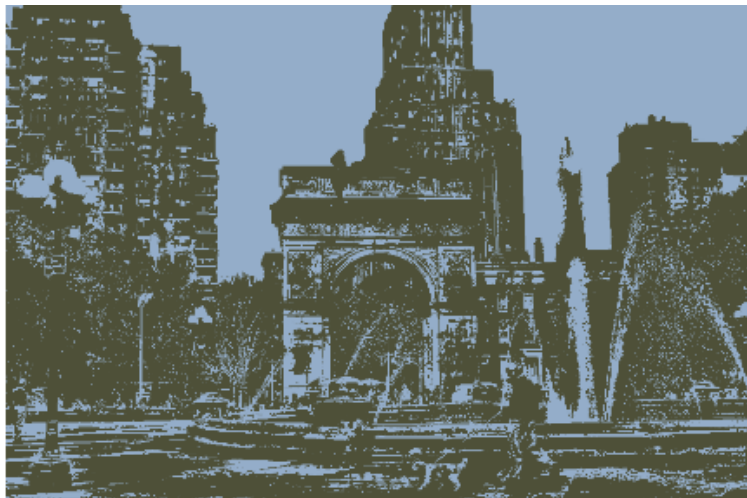
Examples

3 colors



Examples

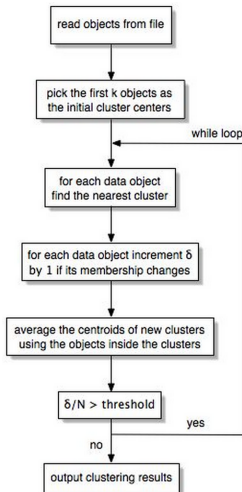
2 colors



Outline

- 1 Introduction
 - The K-Means Algorithm
 - Examples
- 2 **Parallel K-Means**
 - MPI Communication Overview
 - Initialization
- 3 Performance
- 4 Concluding Remarks

MPI Communication Overview



N : number of data objects

K : number of clusters

$objects[N]$: array of data objects

$clusters[K]$: array of cluster centers

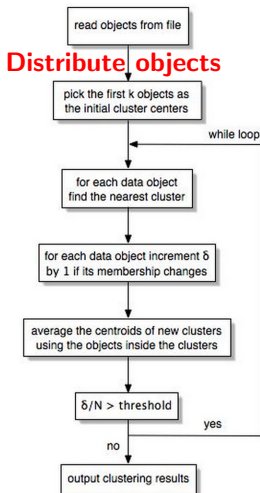
$membership[N]$: array of object memberships

kmeans_clustering()

```

1  while  $\delta/N > \text{threshold}$ 
2     $\delta \leftarrow 0$ 
3    for  $i \leftarrow 0$  to  $N-1$ 
4      for  $j \leftarrow 0$  to  $K-1$ 
5         $\text{distance} \leftarrow |objects[i] - clusters[j]|$ 
6        if  $\text{distance} < d_{\min}$ 
7           $d_{\min} \leftarrow \text{distance}$ 
8           $n \leftarrow j$ 
9        if  $membership[i] \neq n$ 
10          $\delta \leftarrow \delta + 1$ 
11          $membership[i] \leftarrow n$ 
12          $new\_clusters[n] \leftarrow new\_clusters[n] + objects[i]$ 
13          $new\_cluster\_size[n] \leftarrow new\_cluster\_size[n] + 1$ 
14     for  $j \leftarrow 0$  to  $K-1$ 
15        $clusters[j][*] \leftarrow new\_clusters[j][*] / new\_cluster\_size[j]$ 
16        $new\_clusters[j][*] \leftarrow 0$ 
17        $new\_cluster\_size[j] \leftarrow 0$ 
  
```

MPI Communication Overview



N : number of data objects

K : number of clusters

$objects[N]$: array of data objects

$clusters[K]$: array of cluster centers

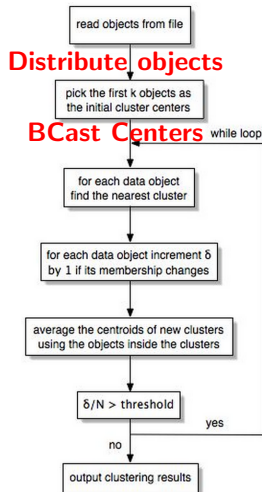
$membership[N]$: array of object memberships

kmeans_clustering()

```

1  while  $\delta/N > \text{threshold}$ 
2     $\delta \leftarrow 0$ 
3    for  $i \leftarrow 0$  to  $N-1$ 
4      for  $j \leftarrow 0$  to  $K-1$ 
5         $\text{distance} \leftarrow |objects[i] - clusters[j]|$ 
6        if  $\text{distance} < d_{\min}$ 
7           $d_{\min} \leftarrow \text{distance}$ 
8           $n \leftarrow j$ 
9        if  $membership[i] \neq n$ 
10          $\delta \leftarrow \delta + 1$ 
11          $membership[i] \leftarrow n$ 
12          $new\_clusters[n] \leftarrow new\_clusters[n] + objects[i]$ 
13          $new\_cluster\_size[n] \leftarrow new\_cluster\_size[n] + 1$ 
14     for  $j \leftarrow 0$  to  $K-1$ 
15        $clusters[j][*] \leftarrow new\_clusters[j][*] / new\_cluster\_size[j]$ 
16        $new\_clusters[j][*] \leftarrow 0$ 
17        $new\_cluster\_size[j] \leftarrow 0$ 
  
```


MPI Communication Overview



N : number of data objects

K : number of clusters

$\text{objects}[N]$: array of data objects

$\text{clusters}[K]$: array of cluster centers

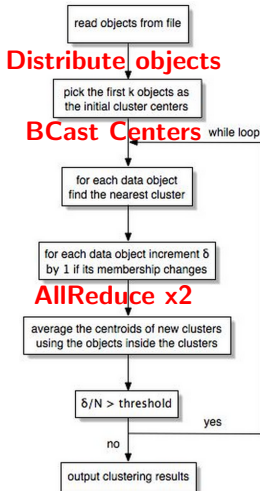
$\text{membership}[N]$: array of object memberships

kmeans_clustering()

```

1  while  $\delta/N > \text{threshold}$ 
2     $\delta \leftarrow 0$ 
3    for  $i \leftarrow 0$  to  $N-1$ 
4      for  $j \leftarrow 0$  to  $K-1$ 
5         $\text{distance} \leftarrow | \text{objects}[i] - \text{clusters}[j] |$ 
6        if  $\text{distance} < d_{\min}$ 
7           $d_{\min} \leftarrow \text{distance}$ 
8           $n \leftarrow j$ 
9        if  $\text{membership}[i] \neq n$ 
10          $\delta \leftarrow \delta + 1$ 
11          $\text{membership}[i] \leftarrow n$ 
12          $\text{new\_clusters}[n] \leftarrow \text{new\_clusters}[n] + \text{objects}[i]$ 
13          $\text{new\_cluster\_size}[n] \leftarrow \text{new\_cluster\_size}[n] + 1$ 
14     for  $j \leftarrow 0$  to  $K-1$ 
15        $\text{clusters}[j][*] \leftarrow \text{new\_clusters}[j][*] / \text{new\_cluster\_size}[j]$ 
16        $\text{new\_clusters}[j][*] \leftarrow 0$ 
17        $\text{new\_cluster\_size}[j] \leftarrow 0$ 
  
```

MPI Communication Overview



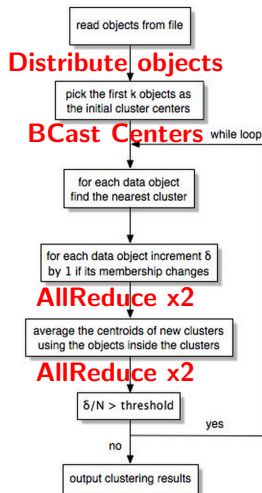
N : number of data objects
 K : number of clusters

$objects[N]$: array of data objects
 $clusters[K]$: array of cluster centers
 $membership[N]$: array of object memberships

```

kmeans_clustering()
1  while  $\delta/N > \text{threshold}$ 
2     $\delta \leftarrow 0$ 
3    for  $i \leftarrow 0$  to  $N-1$ 
4      for  $j \leftarrow 0$  to  $K-1$ 
5         $\text{distance} \leftarrow |objects[i] - clusters[j]|$ 
6        if  $\text{distance} < d_{\min}$ 
7           $d_{\min} \leftarrow \text{distance}$ 
8           $n \leftarrow j$ 
9        if  $membership[i] \neq n$ 
10          $\delta \leftarrow \delta + 1$ 
11          $membership[i] \leftarrow n$ 
12          $new\_clusters[n] \leftarrow new\_clusters[n] + objects[i]$ 
13          $new\_cluster\_size[n] \leftarrow new\_cluster\_size[n] + 1$ 
14     for  $j \leftarrow 0$  to  $K-1$ 
15        $clusters[j][*] \leftarrow new\_clusters[j][*] / new\_cluster\_size[j]$ 
16        $new\_clusters[j][*] \leftarrow 0$ 
17        $new\_cluster\_size[j] \leftarrow 0$ 
  
```

MPI Communication Overview



N : number of data objects

K : number of clusters

$\text{objects}[N]$: array of data objects

$\text{clusters}[K]$: array of cluster centers

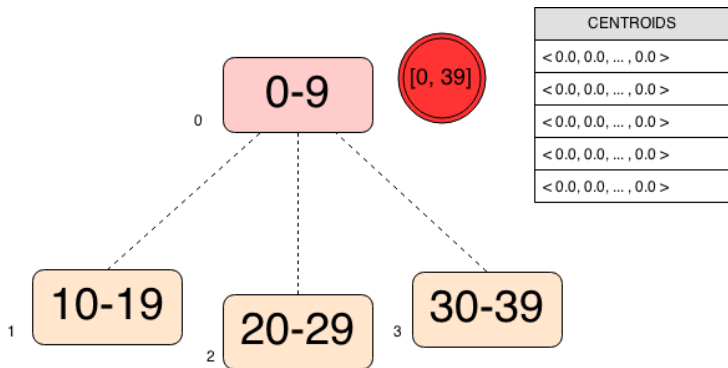
$\text{membership}[N]$: array of object memberships

kmeans_clustering()

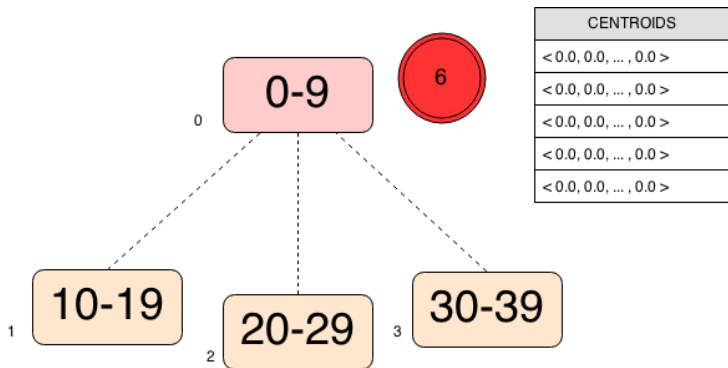
```

1  while  $\delta/N > \text{threshold}$ 
2     $\delta \leftarrow 0$ 
3    for  $i \leftarrow 0$  to  $N-1$ 
4      for  $j \leftarrow 0$  to  $K-1$ 
5         $\text{distance} \leftarrow | \text{objects}[i] - \text{clusters}[j] |$ 
6        if  $\text{distance} < d_{\min}$ 
7           $d_{\min} \leftarrow \text{distance}$ 
8           $n \leftarrow j$ 
9        if  $\text{membership}[i] \neq n$ 
10          $\delta \leftarrow \delta + 1$ 
11          $\text{membership}[i] \leftarrow n$ 
12          $\text{new\_clusters}[n] \leftarrow \text{new\_clusters}[n] + \text{objects}[i]$ 
13          $\text{new\_cluster\_size}[n] \leftarrow \text{new\_cluster\_size}[n] + 1$ 
14     for  $j \leftarrow 0$  to  $K-1$ 
15        $\text{clusters}[j][*] \leftarrow \text{new\_clusters}[j][*] / \text{new\_cluster\_size}[j]$ 
16        $\text{new\_clusters}[j][*] \leftarrow 0$ 
17        $\text{new\_cluster\_size}[j] \leftarrow 0$ 
  
```

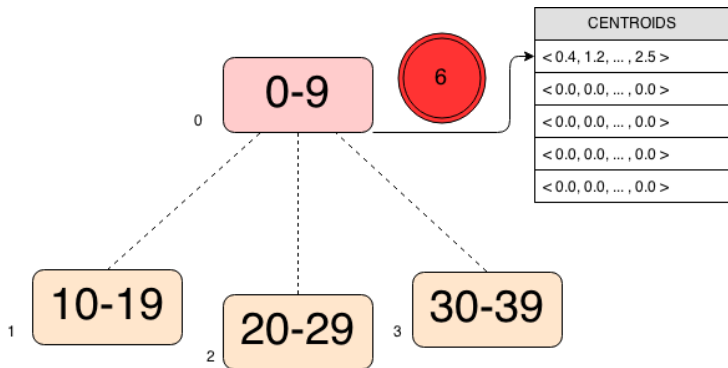
Initialization



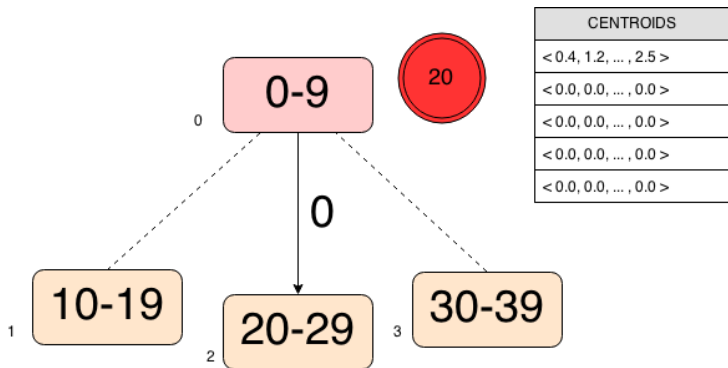
Initialization



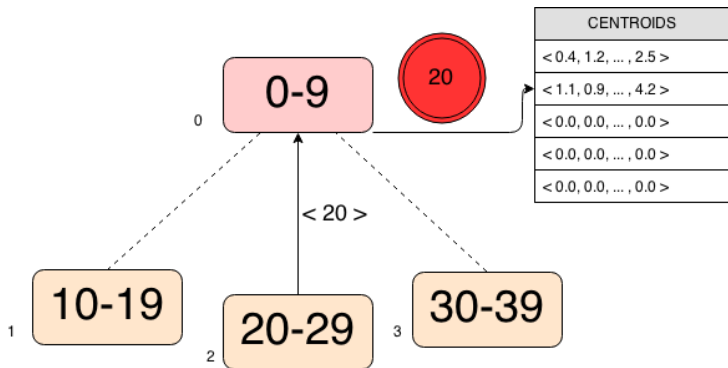
Initialization



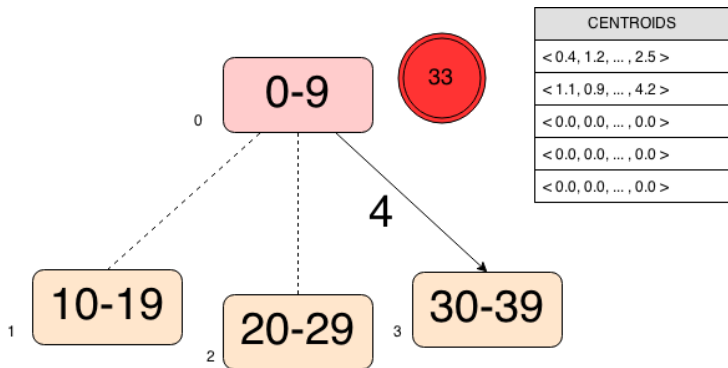
Initialization



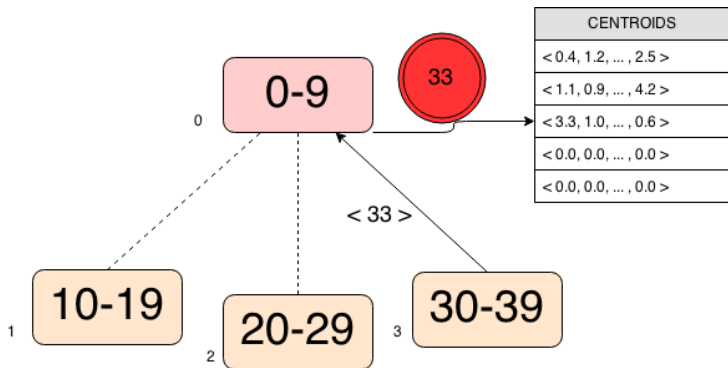
Initialization



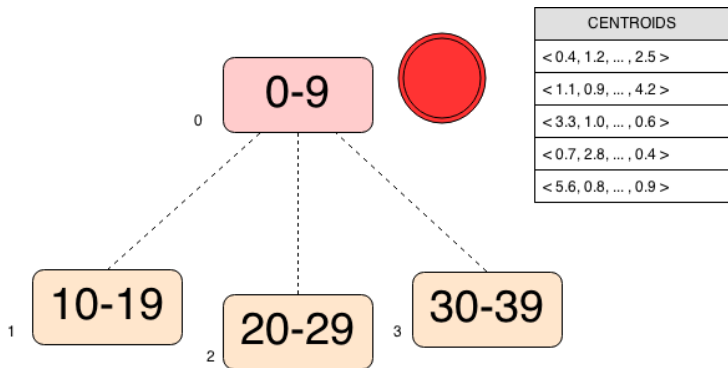
Initialization



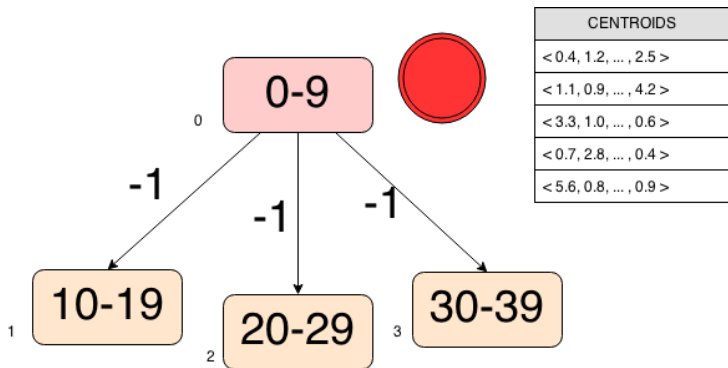
Initialization



Initialization



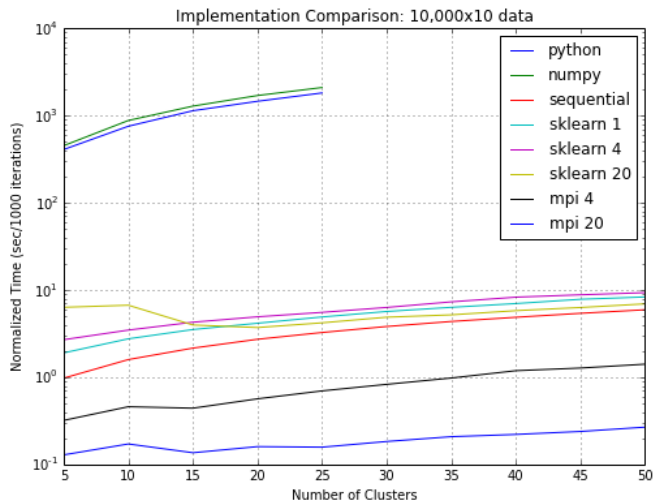
Initialization



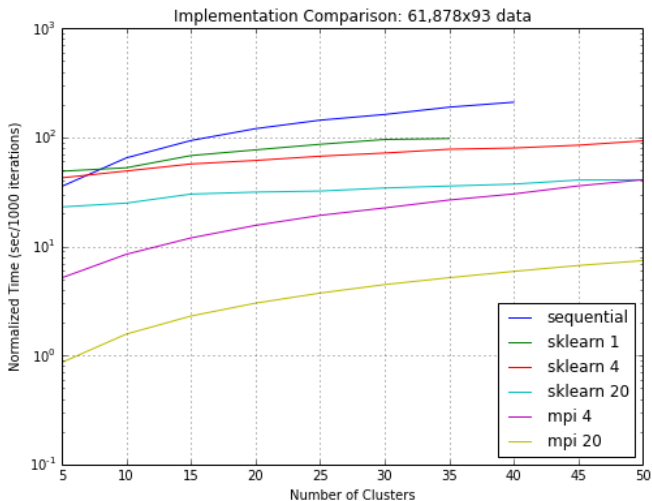
Outline

- 1 Introduction
 - The K-Means Algorithm
 - Examples
- 2 Parallel K-Means
 - MPI Communication Overview
 - Initialization
- 3 Performance
- 4 Concluding Remarks

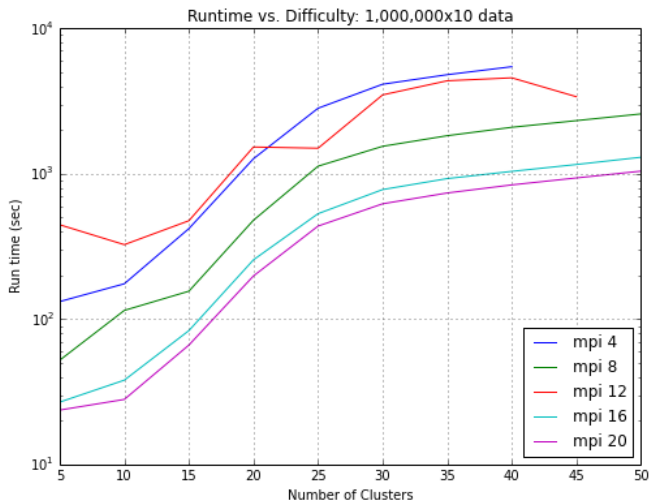
Multiple Implementations



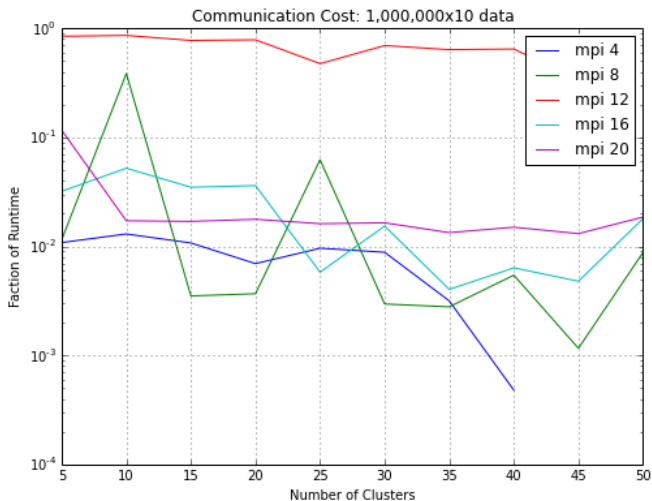
Larger Dataset



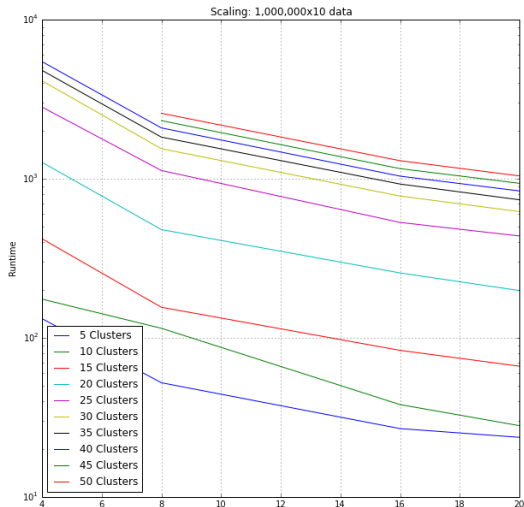
Affect of the Number of Clusters



Communication Cost



Strong Scaling



Outline

- 1 Introduction
 - The K-Means Algorithm
 - Examples
- 2 Parallel K-Means
 - MPI Communication Overview
 - Initialization
- 3 Performance
- 4 Concluding Remarks

Comments

- The best method for parallelization will depend on your data
- There are much better methods for initializing the centroids to converge faster, `kmeans++`
- Reading plain text files with MPI is annoying
- Every now and then communication would hang (locally) for a long time before
- Saw some errors/warnings when running on Mercer for larger data that I didn't see previously

huhh??

```
Computation time: 10412.692955s
Communication time: 68.016616s
```

```
-----
_orterrun noticed that process rank 6 with PID 10686 on node compute-15-9 exited on signal 11 (Segmentation fault).
-----
```

```
WARNING: a request was made to bind a process. While the system
supports binding the process itself, at least one node does NOT
support binding memory to the process location.
```

```
Node: compute-15-9
```

```
This usually is due to not having the required NUMA support installed
on the node. In some Linux distributions, the required support is
contained in the libnumactl and libnumactl-devel packages.
This is a warning only; your job will continue, though performance may be degraded.
```

```
-----
The library attempted to open the following supporting CUDA libraries,
but each of them failed.  CUDA-aware support is disabled.
libcuda.so.1: cannot open shared object file: No such file or directory
/usr/lib64/libcuda.so.1: cannot open shared object file: No such file or directory
If you are not interested in CUDA-aware support, then run with
--mca mpi_cuda_support 0 to suppress this message.  If you are interested
in CUDA-aware support, then try setting LD_LIBRARY_PATH to the location
of libcuda.so.1 to get passed this issue.
```

```
-----
[END]
```

Further Reading I



Arthur, David and Manthey, Bodo and Roglin, H

k-Means has polynomial smoothed complexity

Foundations of Computer Science, 2009. FOCS'09. 50th Annual IEEE Symposium on. (2009) 405-414



MacQueen, James and others

Some methods for classification and analysis of multivariate observations

Proceedings of the fifth Berkeley symposium on mathematical statistics and probability. (1967) 281-297



Arthur, David and Vassilvitskii, Sergei

k-means++: The advantages of careful seeding

Proceedings of the eighteenth annual ACM-SIAM symposium on Discrete algorithms. (2007) 1027-1035