



K-Means Clustering

Using Map-Reduce

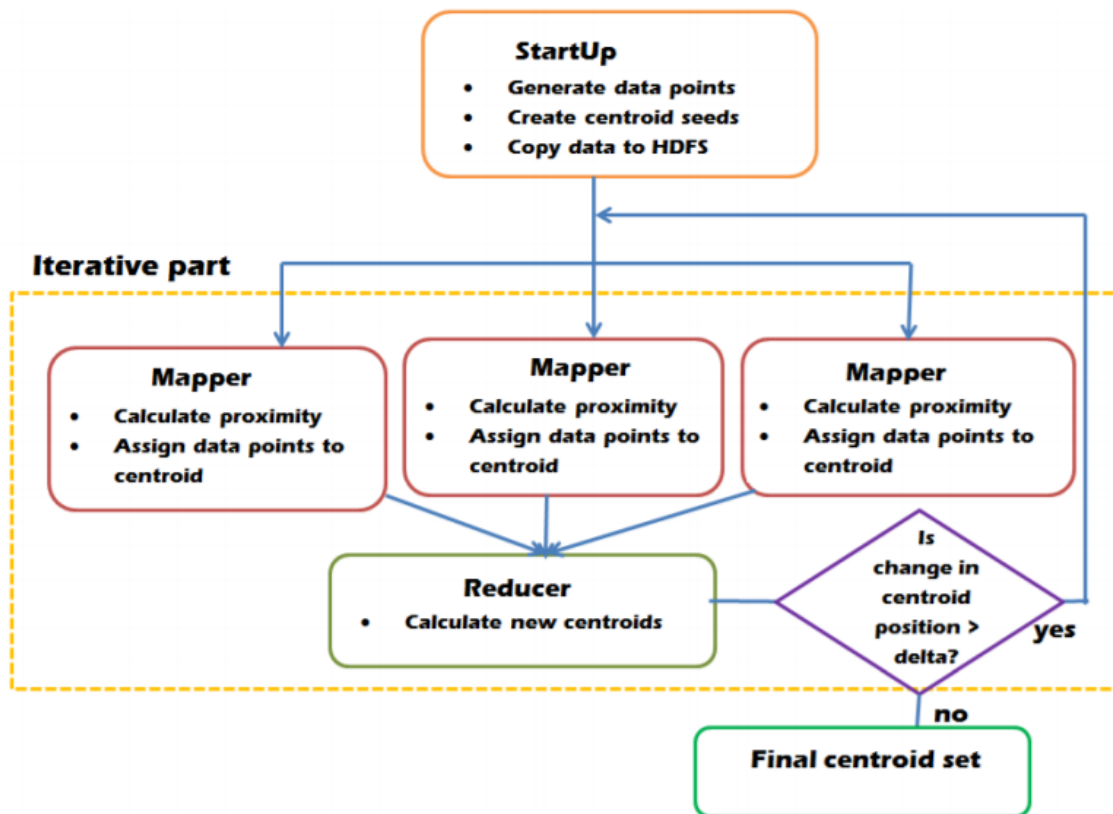
Team Members:

Rahul Sathyajit
Karthick Raja
Akshar Vashist

Contents

Structure of the program
Algorithm for Startup
Algorithm for Mapper
Algorithm for Reducer
Algorithm for Simple Assignment
Algorithm for Centroid Calculator

Structure of the Algorithm



Algorithm for Startup

Require:

- A set of d-dimensional objects $X = \{x_1, x_2, \dots, x_n\}$
- k-number of clusters where $k < n$
- initial set of centroids $C = \{c_1, c_2, \dots, c_k\}$
- δ convergence delta

Output: A new set of centroids, number of iterations, final clusters and the time taken to converge

```
1: current_centroids  $\leftarrow C$ 
2: initialise numIter, timetoConverge, finalClusters
3: startTime  $\leftarrow$  currentTime()
4:  $C' \leftarrow$  call SimpleAssignment
5: new_centroids  $\leftarrow C'$ 
6: numIter  $\leftarrow 1$ 
7: while change(new_centroids, current_centroids)  $> \delta$  do
8:   current_centroids  $\leftarrow$  new_centroids
9:    $C' \leftarrow$  call SimpleAssignment
10:  numIter  $\leftarrow$  numIter + 1
11:  current_centroids  $\leftarrow C'$ 
12: end while
13: endTime  $\leftarrow$  currentTime()
14: timetoConverge  $\leftarrow$  (endTime – startTime)
15: perform outlierRemoval
16: finalClusters  $\leftarrow$  perform finalClustering
17: writeStatistics
18: return current_centroids , numIter, finalClusters, timetoConverge
```

Algorithm for Simple Assignment

Require:

- A set of d-dimensional objects of $\{x_1, x_2, \dots, x_n\}$
- initial set of centroids $C = \{c_1, c_2, \dots, c_k\}$

Output: A dictionary list consisting of each centroid and the objects assigned to them.
This list is passed as an argument to the CentroidCalculator program.

```
1:  $M_1 \leftarrow \{x_1, x_2, \dots, x_n\}$ 
2:  $current\_centroids \leftarrow C$ 
3:  $distance(p, q) = \sqrt{\sum_{i=1}^d (p_i - q_i)^2}$  where  $p_i$  (or  $q_i$ ) is the coordinate of p (or q) in dimension
   i  $v \leftarrow \{\}$ 
4: for all  $x_i \in M_1$  such that  $1 \leq i \leq n$  do
5:    $bestCentroid \leftarrow null$ 
6:    $minDist \leftarrow \infty$ 
7:   for all  $c$  in  $current\_centroids$  do
8:      $dist \leftarrow distance(x_i, c)$ 
9:     if  $bestCentroid = null \parallel dist < minDist$  then
10:       $bestCentroid \leftarrow c$ 
11:       $minDist \leftarrow dist$ 
12:     end if
13:   end for
14:    $v[bestCentroid] \leftarrow (x_i)$ 
15:    $i += 1$ 
16: end for
17: call  $CentroidCalculator(v)$ 
```

Algorithm for centroid calculator

Require:

Input: List of centroids and their assigned objects

Output: List of new centroids

```
1:  $v = \{\}$ 
2:  $v = \text{outputlist}$  from SimpleAssignment
3:  $\text{newCentroidList} \leftarrow \text{null}$ 
4: for all  $\text{centroid} \in v$  do
5:    $\text{newCentroid}, \text{sumofObjects}, \text{numofObjects} \leftarrow \text{null}$ 
6:   for all  $\text{object} \in v[\text{centroid}]$  do
7:      $\text{sumofObjects} += \text{object}$ 
8:      $\text{numofObjects} += 1$ 
9:   end for
10:   $\text{newCentroid} \leftarrow (\text{sumofObjects} \div \text{numofObjects})$ 
11:   $\text{newCentroidList} << (\text{newCentroid})$ 
12: end for
13: return  $\text{newCentroidList}$ 
```

Algorithm for Mapper

Require:

- A subset of d-dimensional objects of $\{x_1, x_2, \dots, x_n\}$ in each mapper
- initial set of centroids $C = \{c_1, c_2, \dots, c_k\}$

Output: A list of centroids and objects assigned to each centroid separated by tab. This list is written locally one line per data point and read by the Reducer program.

```
1:  $M_1 \Leftarrow \{x_1, x_2, \dots, x_m\}$ 
2:  $current\_centroids \Leftarrow C$ 
3:  $distance(p, q) = \sqrt{\sum_{i=1}^d (p_i - q_i)^2}$  where  $p_i$  (or  $q_i$ ) is the coordinate of p (or q) in dimension
   i
4: for all  $x_i \in M_1$  such that  $1 \leq i \leq m$  do
5:    $bestCentroid \Leftarrow null$ 
6:    $minDist \Leftarrow \infty$ 
7:   for all  $c \in current\_centroids$  do
8:      $dist \Leftarrow distance(x_i, c)$ 
9:     if  $bestCentroid = null \parallel dist < minDist$  then
10:       $minDist \Leftarrow dist$ 
11:       $bestCentroid \Leftarrow c$ 
12:     end if
13:   end for
14:    $outputlist \leftarrow (bestCentroid, x_i)$ 
15:    $i += 1$ 
16: end for
17: return  $outputlist$ 
```

Algorithm for Reducer

Require:

Input: (key, value) where key = bestCentroid and value = objects assigned to the centroids by the mapper.

Output: (key, value) where key = oldcentroid and value = newBestCentroid which is the new centroid value calculated for that bestCentroid.

```
1: outputlist  $\leftarrow$  outputlists from mappers
2: v  $\leftarrow$  {}
3: newCentroidList  $\leftarrow$  null
4: for all y  $\in$  outputlist do
5:   centroid  $\leftarrow$  y.key
6:   object  $\leftarrow$  y.value
7:   v[centroid]  $\leftarrow$  object
8: end for
9: for all centroid  $\in$  v do
10:  newCentroid, sumofObjects, numofObjects  $\leftarrow$  null
11:  for all object  $\in$  v[centroid] do
12:    sumofObjects  $+=$  object
13:    numofObjects  $+=$  1
14:  end for
15:  newCentroid  $\leftarrow$  (sumofObjects  $\div$  numofObjects)
16:  newCentroidList  $\ll$  (newCentroid)
17: end for
18: return newCentroidList
```

Reference:

D. Gillick, A. Faria, and J. DeNero. Mapreduce: Distributed computing for machine learning, 2006

W. Zhao, H. Ma, and Q. He. Parallel k-means clustering based on mapreduce. Cloud Computing, pages 674–679, 2009.

K-means Clustering Using Hadoop MapReduce, Grace Nila Ramamoorthy