

ROADMAP

Kmeans

Mapreduce implementation

Demonstration

Results and statistics

KMEANS

Given a set of observations (x_1, x_2, \dots, x_n) , where each observation is a d -dimensional real vector, k -means clustering aims to partition the n observations into k sets ($k \leq n$ and $k > 0$) $S = \{S_1, S_2, \dots, S_k\}$ based on features/attributes of data points so as to minimize the within-cluster sum of squares (WCSS):

$$\arg \min_{\mathbf{S}} \sum_{i=1}^k \sum_{x_j \in S_i} \|x_j - \mu_i\|^2$$

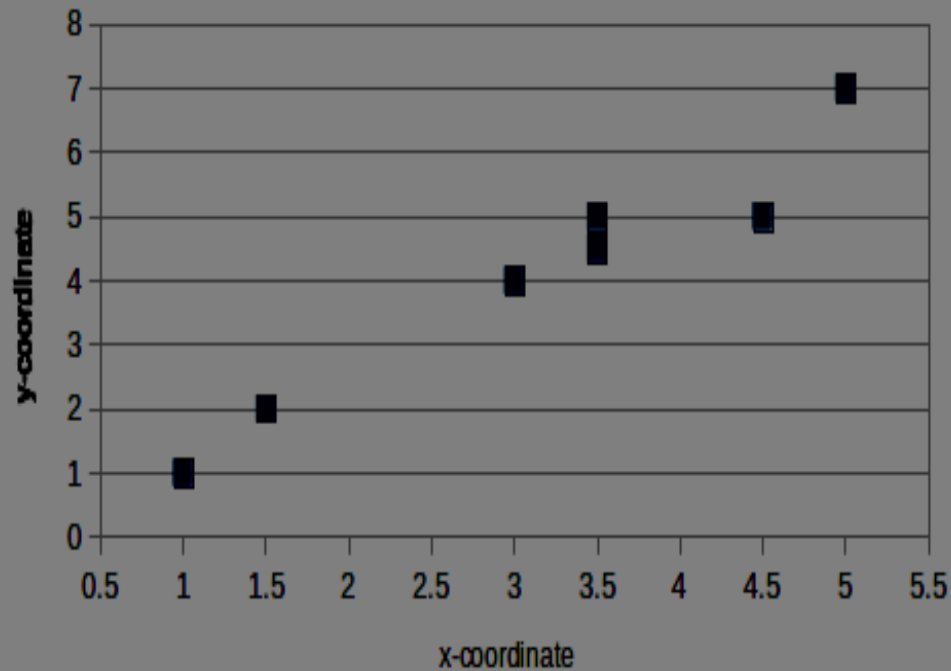
where μ_i is the geometric centroid of data points in S_i and x_j is the vector representing j th data point.

MAPREDUCE IMPLEMENTATION

DEMONSTRATION

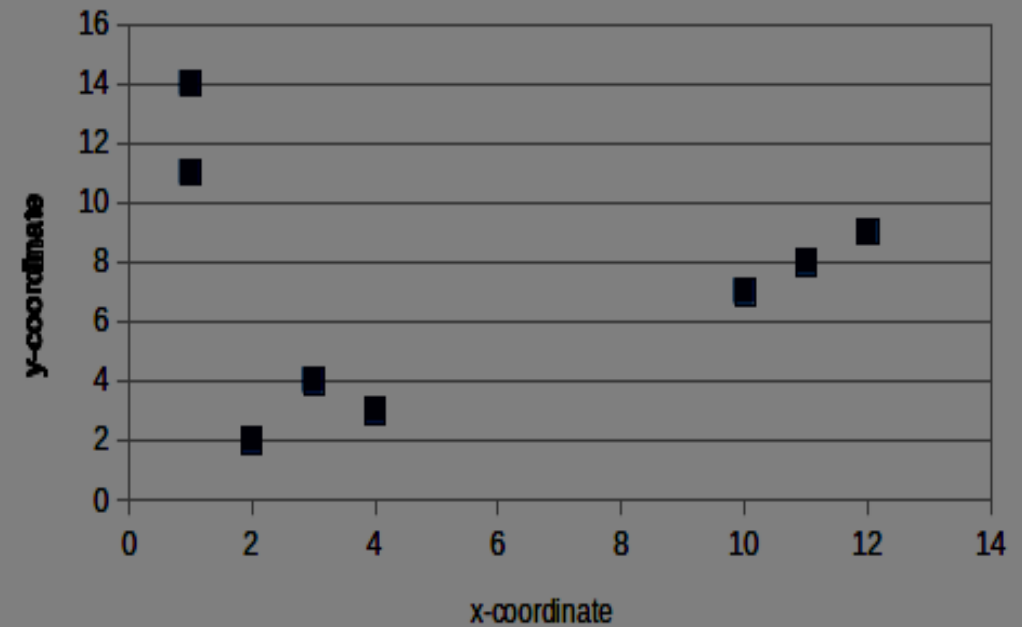
KMeans Clustering

points=7, dimensions=2, clusters=2



KMeans Clustering

points=8, dimensions=2, clusters=3



RESULTS AND STATISTICS

•INPUT DATA SET

- Number of points(n) = 1 million
- Dimensions(d) = 10
- Number of clusters(k) = 10
- Size of dataset = 31.2 MB

•RUN TIME ON LOCAL MACHINE

- For #iterations = 1, runtime = 0:09:33.797533

•RUN TIME ON HADOOP

- For #iterations = 1, runtime = 0:11:45.287343