

CS 412 Introduction to Machine Learning

# K-means

Instructor: Wei Tang

Department of Computer Science  
University of Illinois at Chicago  
Chicago IL 60607

<https://tangw.people.uic.edu>  
[tangw@uic.edu](mailto:tangw@uic.edu)

Slides credit: Xinhua Zhang

# k-Means Clustering

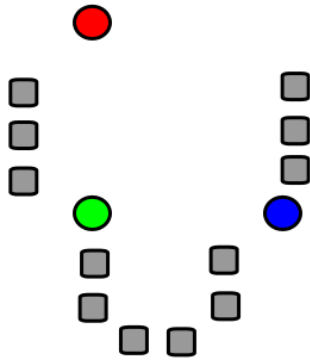
- Find  $k$  reference vectors (prototypes/codebook vectors/codewords) which best represent data
- Sample  $\mathcal{X} = \{\mathbf{x}^t\}_{t=1}^N$ . Reference vectors:  $\mathbf{m}_j$  ( $j = 1, \dots, k$ )
- Use nearest (most similar) reference: code book

$$\|\mathbf{x}^t - \mathbf{m}_i\| = \min_j \|\mathbf{x}^t - \mathbf{m}_j\|$$

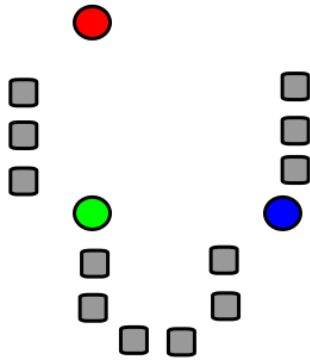
- Reconstruction error  $E(\{\mathbf{m}_i\}_{i=1}^k | \mathcal{X}) = \sum_t \sum_i b_i^t \|\mathbf{x}^t - \mathbf{m}_i\|^2$

no analytic minimizer  
NP-hard to optimize  $\{m_i\}$

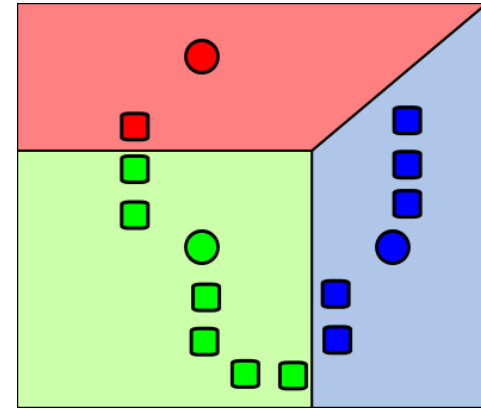
$$b_i^t = \begin{cases} 1 & \text{if } \|\mathbf{x}^t - \mathbf{m}_i\| = \min_j \|\mathbf{x}^t - \mathbf{m}_j\| \\ 0 & \text{otherwise} \end{cases}$$



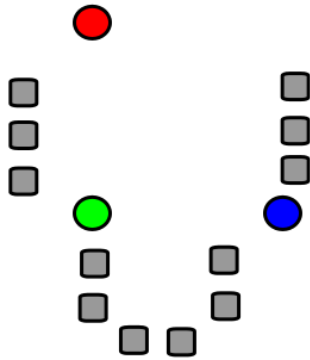
1. Select initial  
centroids at random



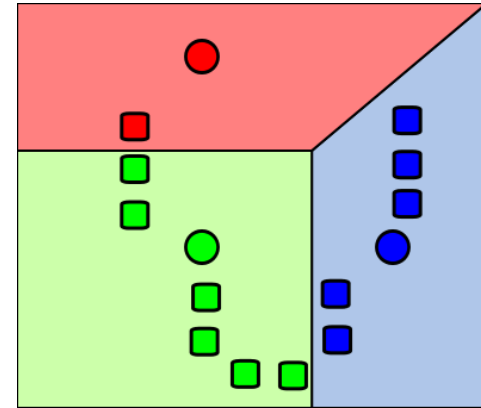
1. Select initial centroids at random



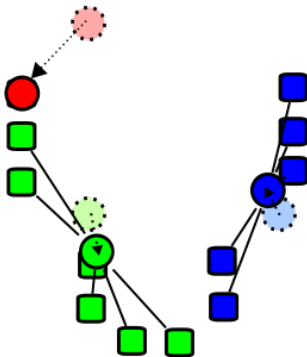
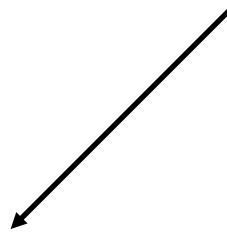
2. Assign each object to the cluster with the nearest centroid.



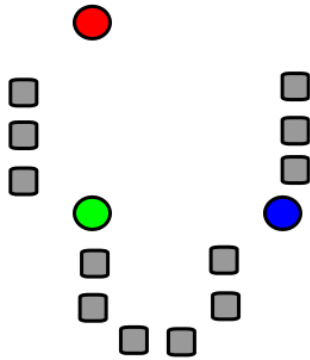
1. Select initial centroids at random



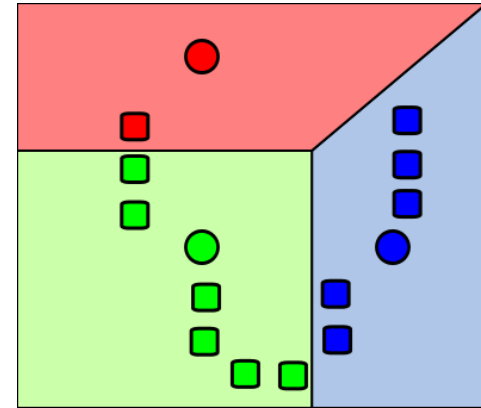
2. Assign each object to the cluster with the nearest centroid.



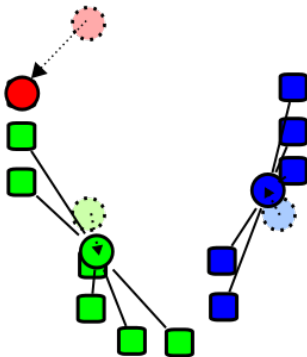
3. Compute each centroid as the mean of the objects assigned to it (go to 2)



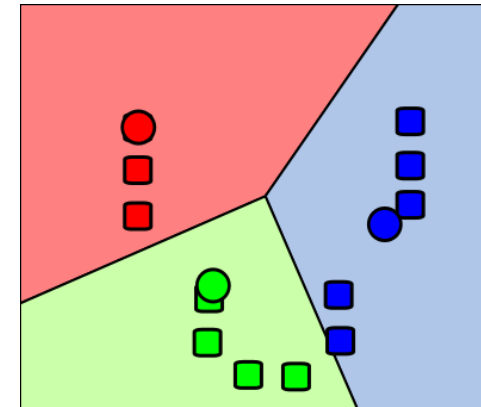
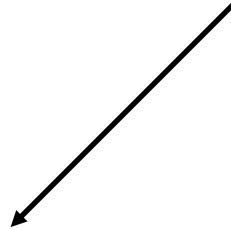
1. Select initial centroids at random



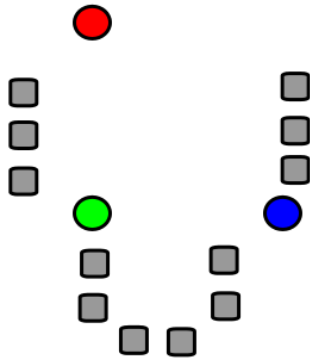
2. Assign each object to the cluster with the nearest centroid.



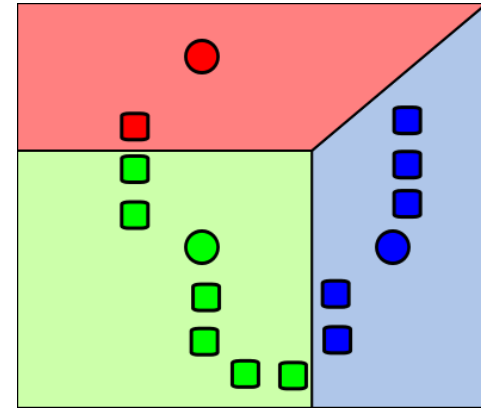
3. Compute each centroid as the mean of the objects assigned to it (go to 2)



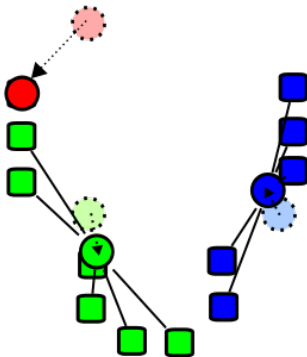
2. Assign each object to the cluster with the nearest centroid.



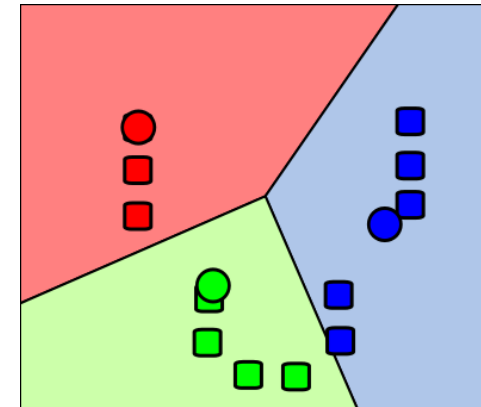
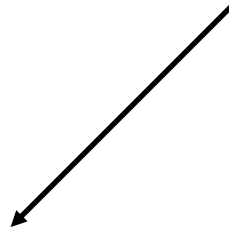
1. Select initial centroids at random



2. Assign each object to the cluster with the nearest centroid.



3. Compute each centroid as the mean of the objects assigned to it (go to 2)



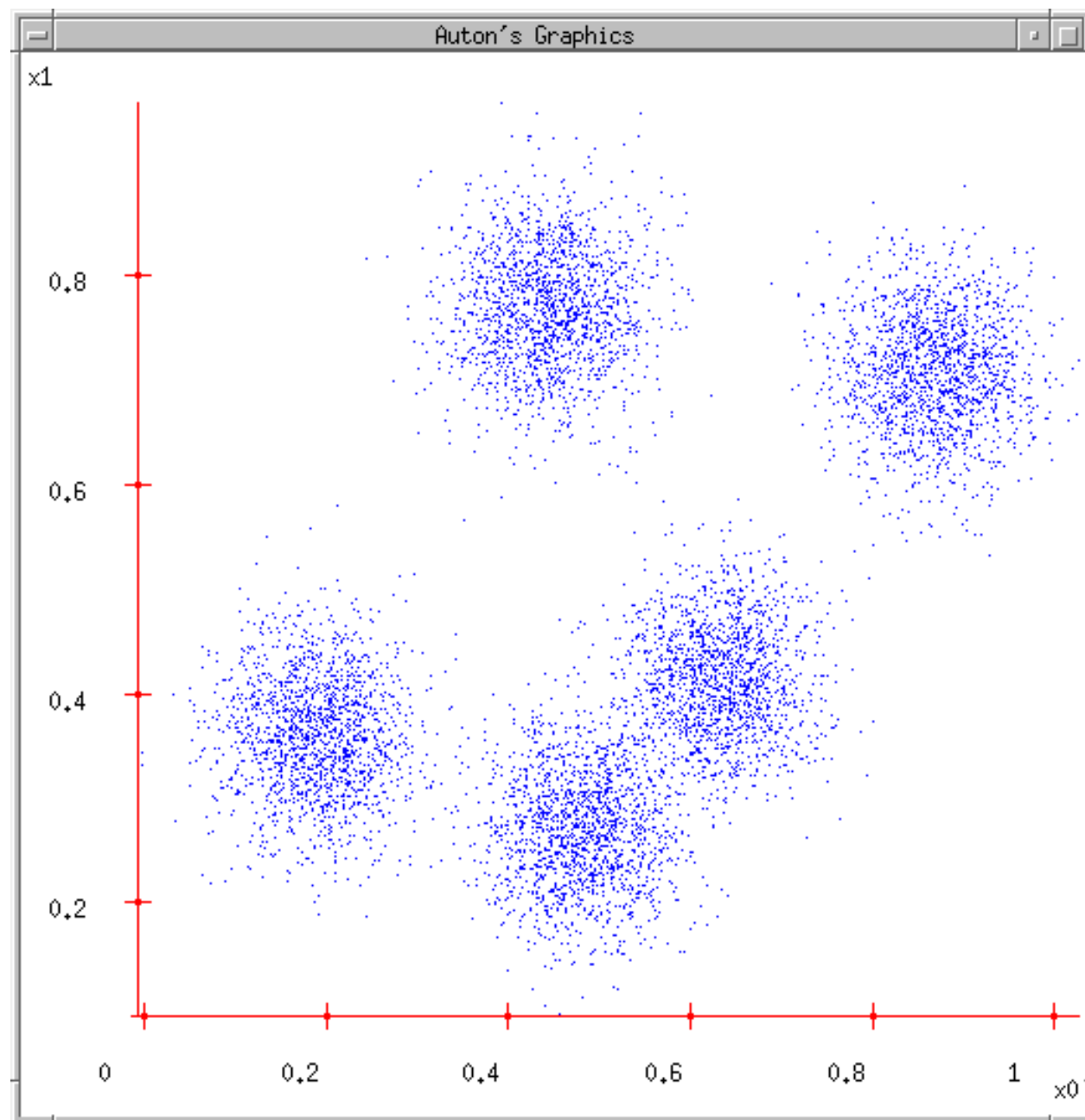
2. Assign each object to the cluster with the nearest centroid.

# K-means Clustering

Given  $k$ :

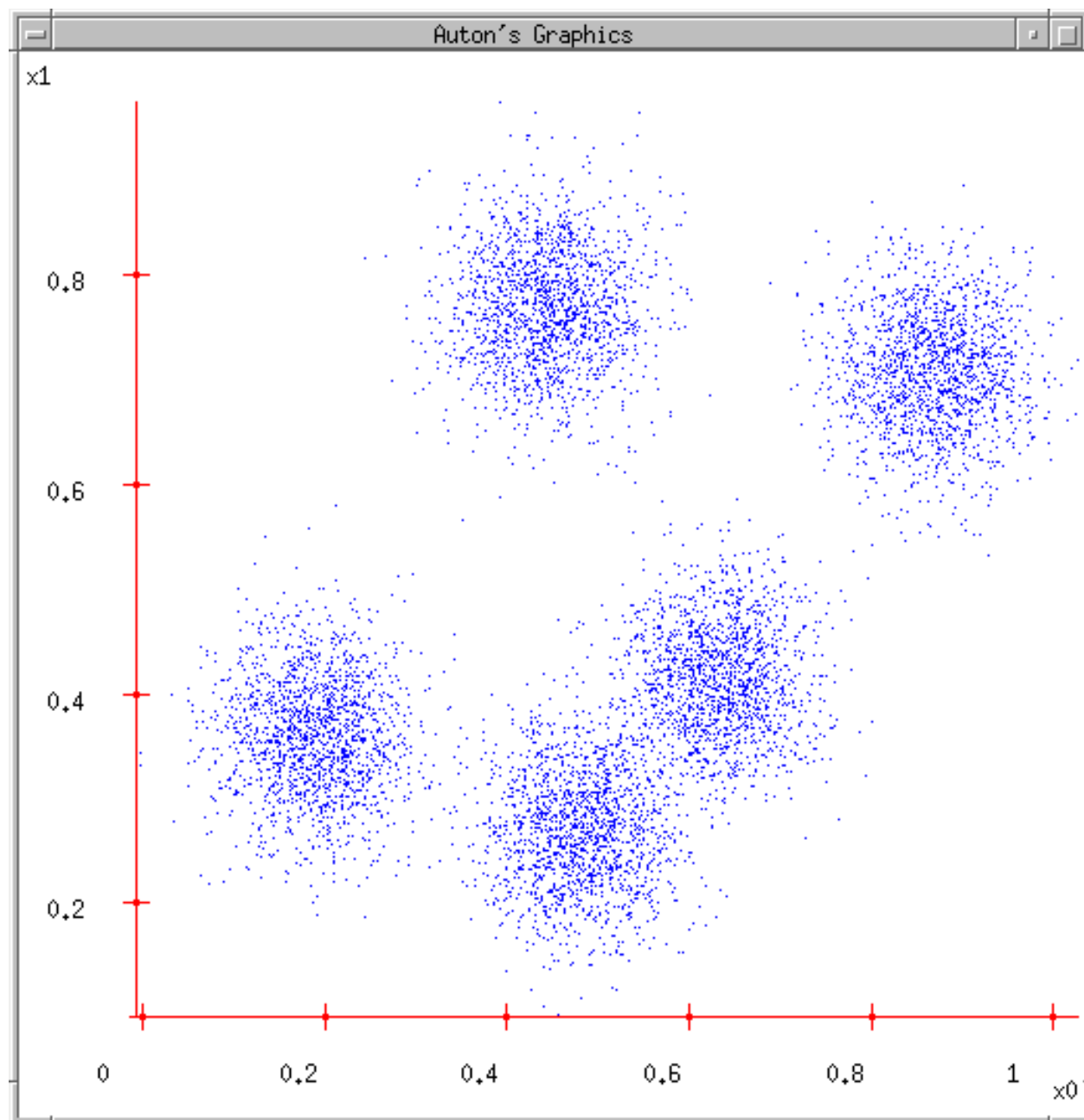
1. Select initial centroids at random.
2. Assign each object to the cluster with the nearest centroid.
3. Compute each centroid as the mean of the objects assigned to it.
4. Repeat previous 2 steps until no change.





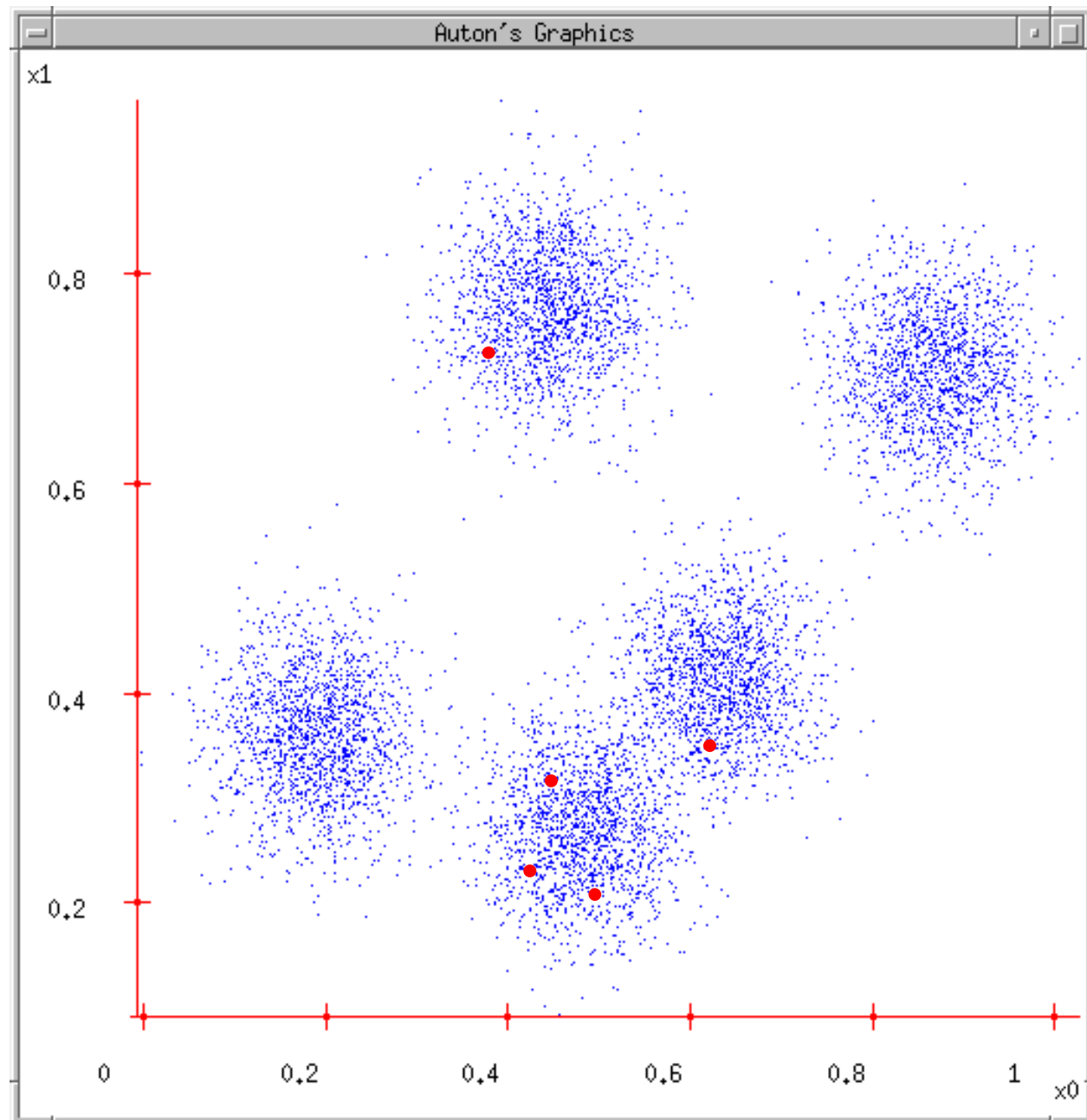
# K-means

1. Ask user how many clusters they'd like.  
(e.g.  $k=5$ )



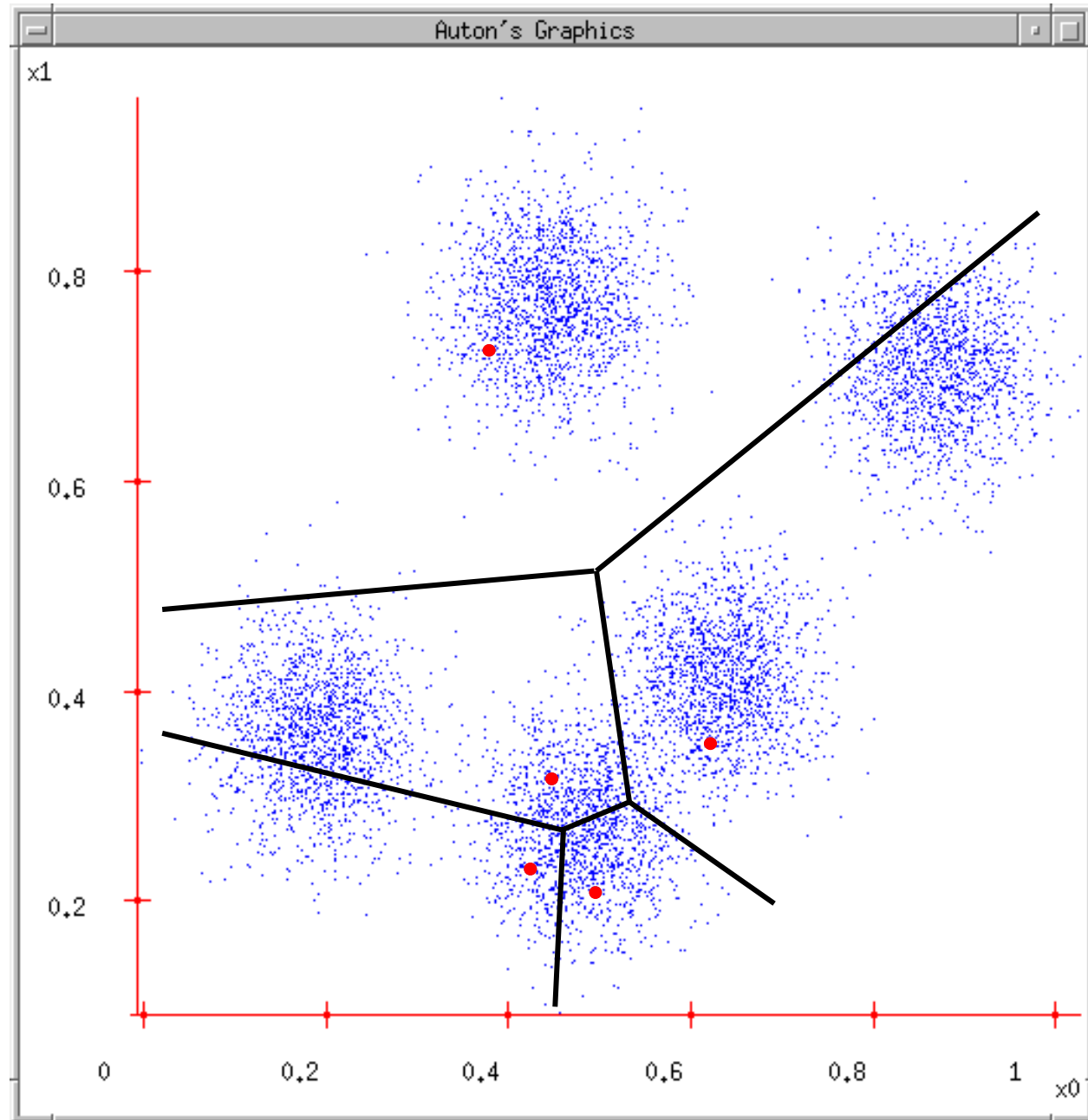
# K-means

1. Ask user how many clusters they'd like.  
(*e.g.  $k=5$* )
2. Randomly guess  $k$  cluster Center locations



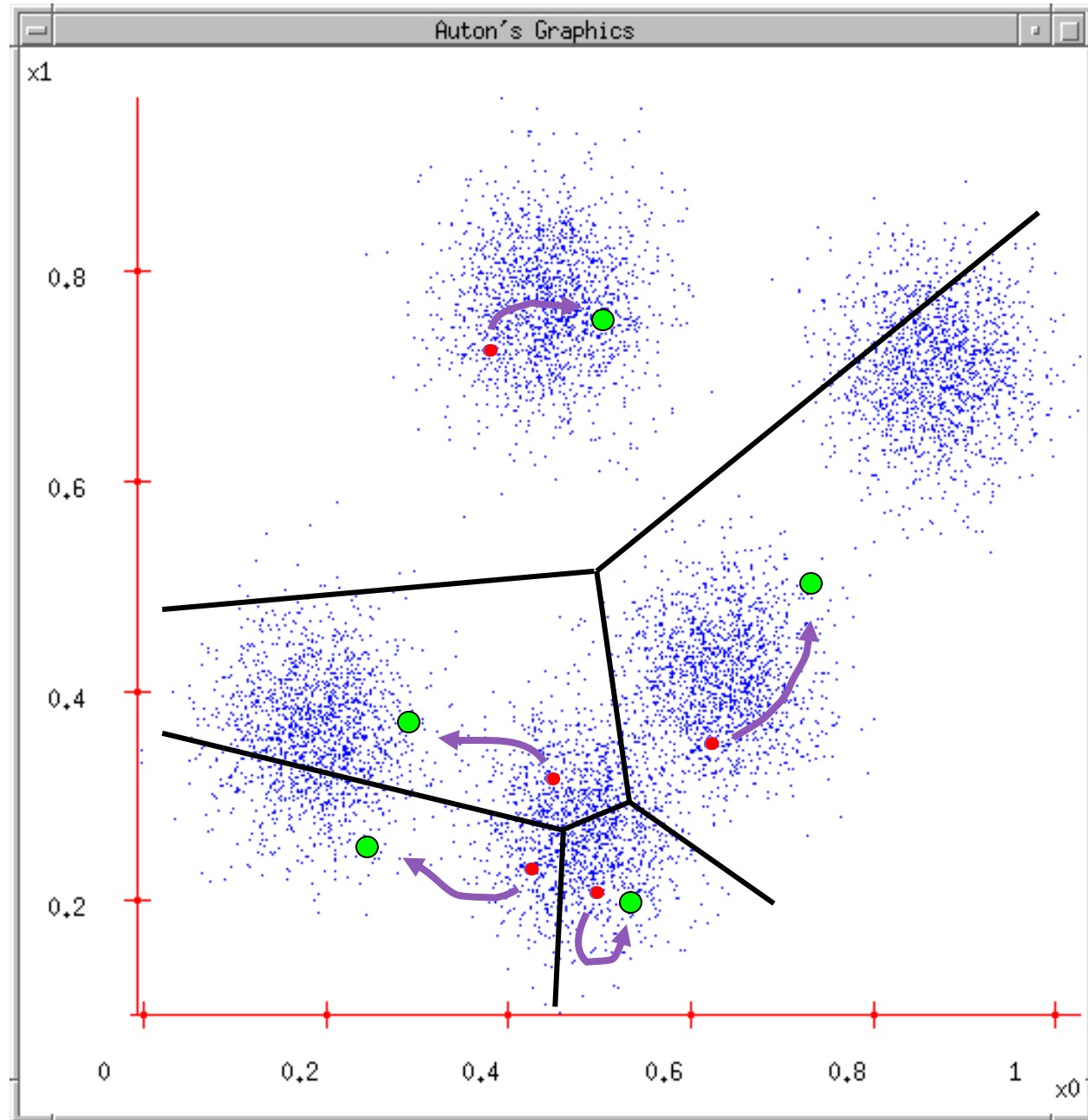
# K-means

1. Ask user how many clusters they'd like.  
(*e.g.  $k=5$* )
2. Randomly guess  $k$  cluster Center locations
3. Each datapoint finds out which Center it's closest to. (Thus each Center "owns" a set of datapoints)



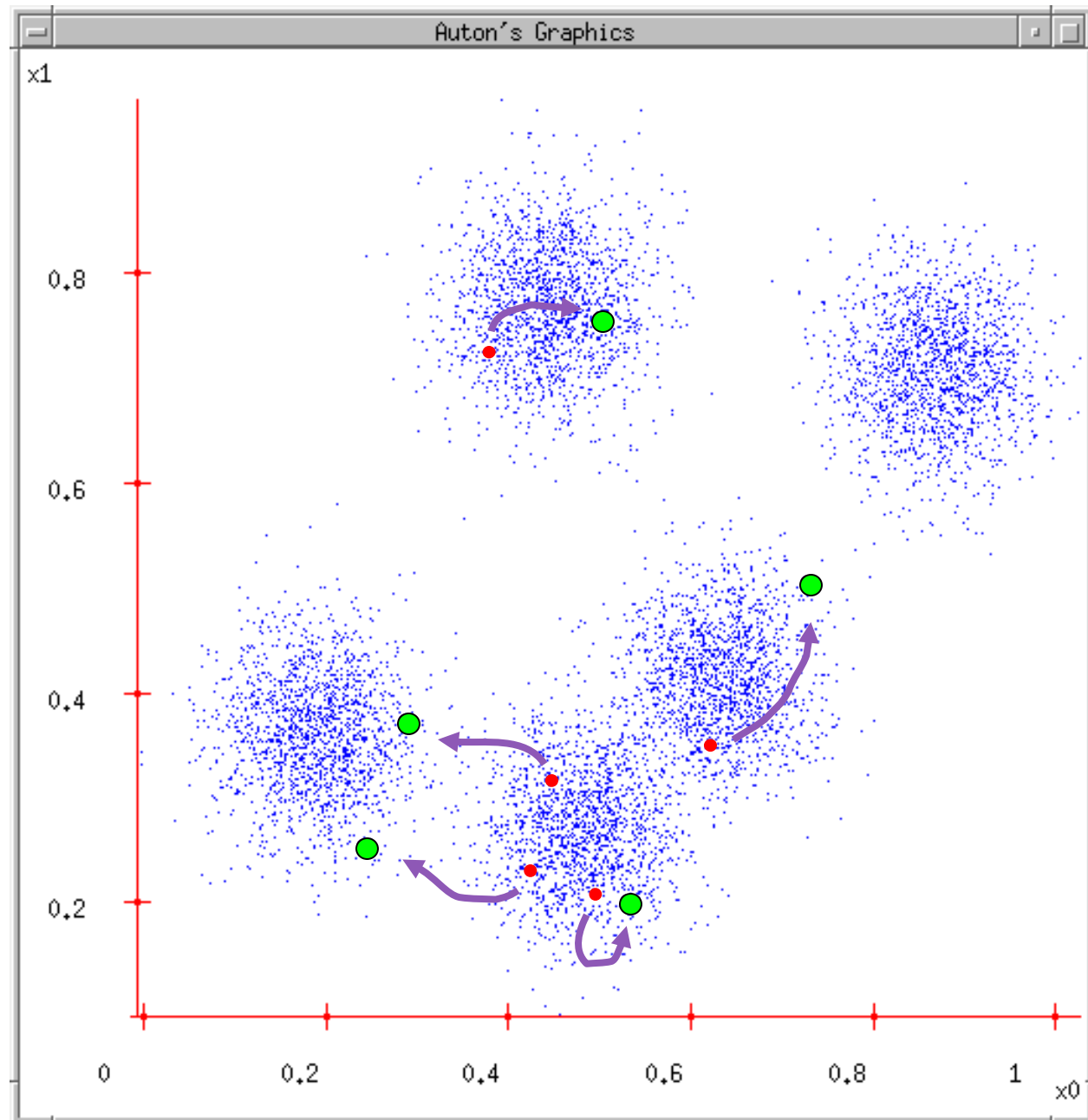
# K-means

1. Ask user how many clusters they'd like.  
(e.g.  $k=5$ )
2. Randomly guess  $k$  cluster Center locations
3. Each datapoint finds out which Center it's closest to.
4. Each Center finds the centroid of the points it owns



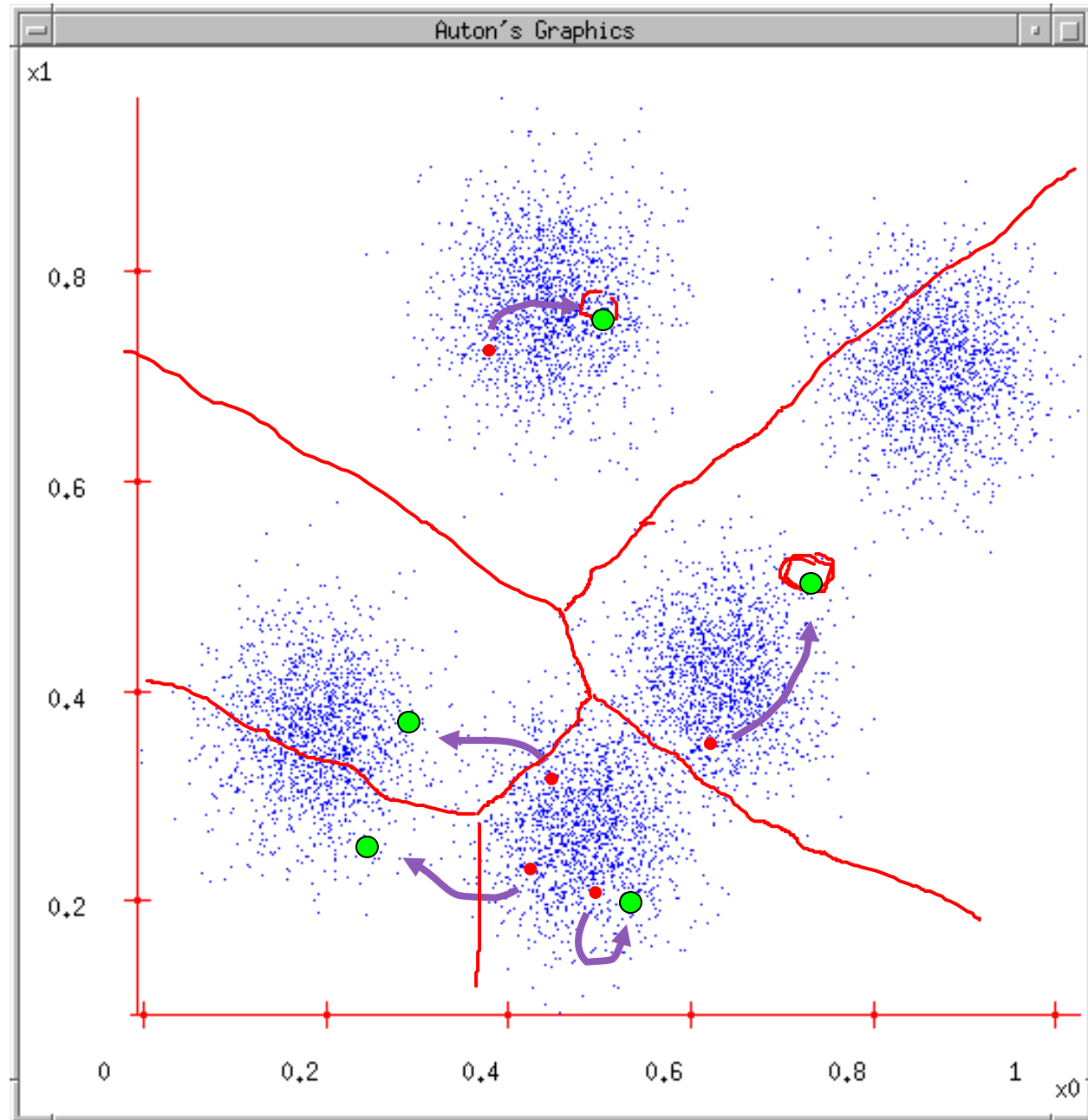
# K-means

1. Ask user how many clusters they'd like.  
(*e.g.  $k=5$* )
2. Randomly guess  $k$  cluster Center locations
3. Each datapoint finds out which Center it's closest to.
4. Each Center finds the centroid of the points it owns...
5. ...and jumps there
6. ...Repeat until terminated!



# K-means

1. Ask user how many clusters they'd like.  
(e.g.  $k=5$ )
2. Randomly guess  $k$  cluster Center locations
3. Each datapoint finds out which Center it's closest to.
4. Each Center finds the centroid of the points it owns...
5. ...and jumps there
6. ...Repeat until terminated!



# *k*-means Clustering

Initialize  $\mathbf{m}_i, i = 1, \dots, k$ , for example, to  $k$  random  $\mathbf{x}^t$

Repeat

For all  $\mathbf{x}^t \in \mathcal{X}$

$$b_i^t \leftarrow \begin{cases} 1 & \text{if } \|\mathbf{x}^t - \mathbf{m}_i\| = \min_j \|\mathbf{x}^t - \mathbf{m}_j\| \\ 0 & \text{otherwise} \end{cases}$$

For all  $\mathbf{m}_i, i = 1, \dots, k$

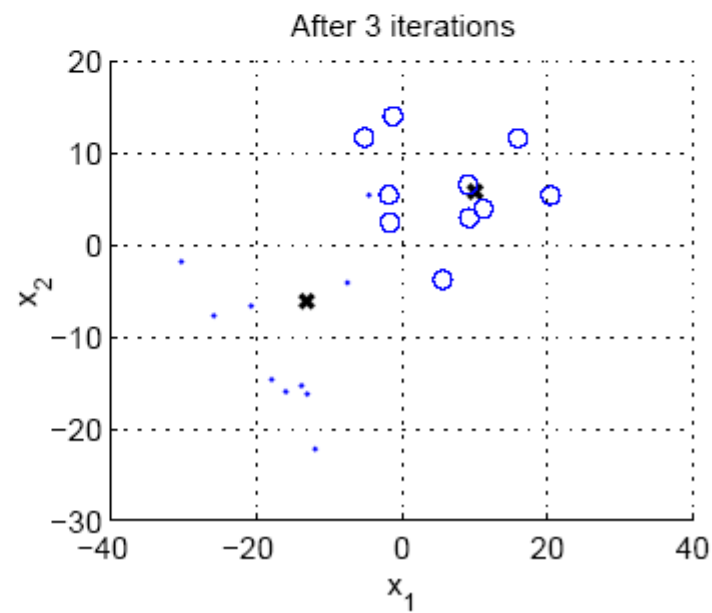
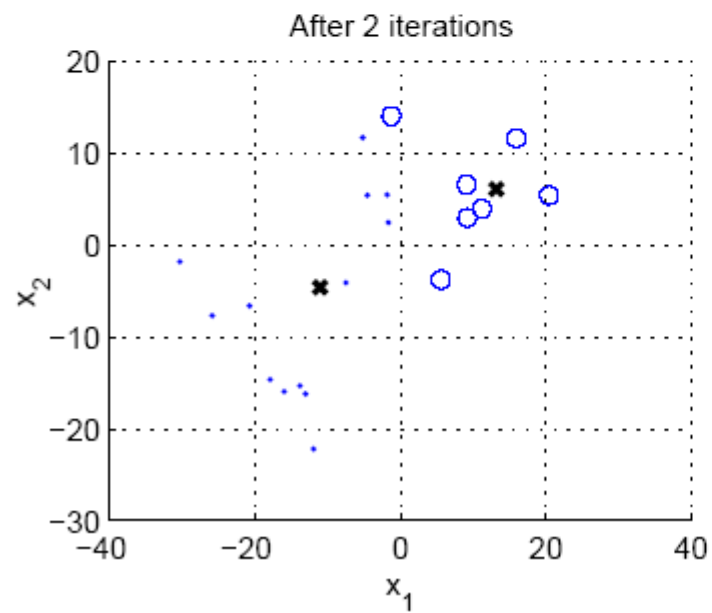
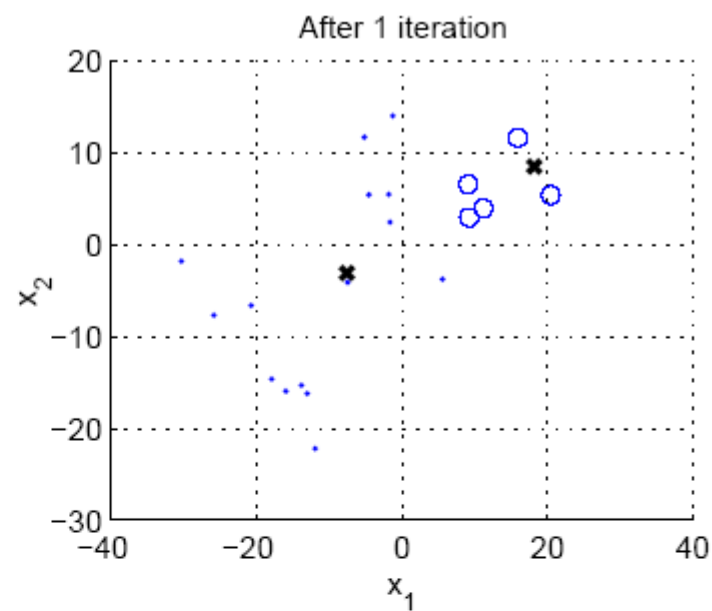
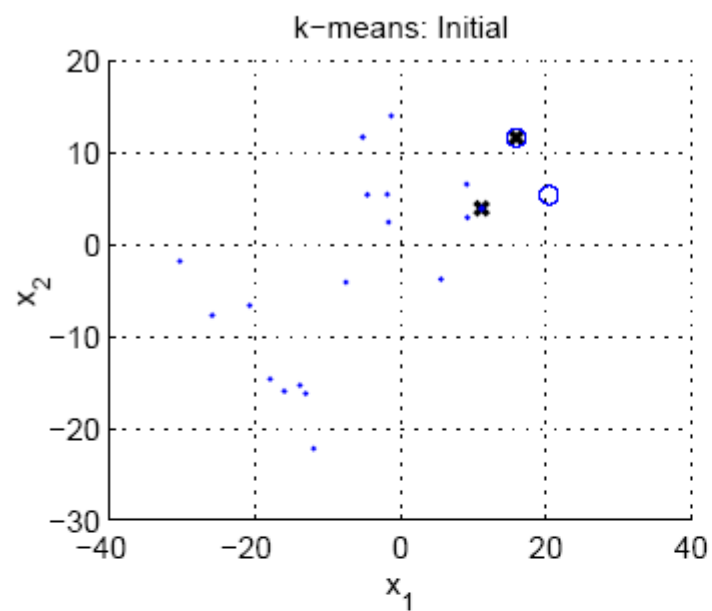
$$\mathbf{m}_i \leftarrow \sum_t b_i^t \mathbf{x}^t / \sum_t b_i^t$$

Until  $\mathbf{m}_i$  converge

$$E(\{\mathbf{m}_i\}_{i=1}^k | \mathcal{X}) = \sum_t \sum_i b_i^t \|\mathbf{x}^t - \mathbf{m}_i\|^2$$

$$b_i^t = \begin{cases} 1 & \text{if } \|\mathbf{x}^t - \mathbf{m}_i\| = \min_j \|\mathbf{x}^t - \mathbf{m}_j\| \\ 0 & \text{otherwise} \end{cases}$$





# Local procedure

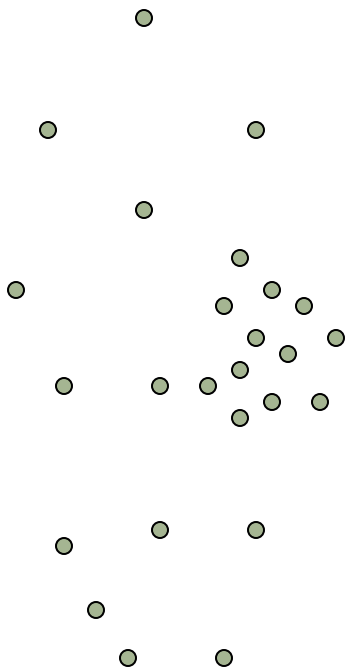
- May converge to suboptimal

$$\min_{\{m_i\}} E(\{\mathbf{m}_i\}_{i=1}^k | \mathcal{X}) = \sum_t \sum_i b_i^t \|\mathbf{x}^t - \mathbf{m}_i\|^2$$

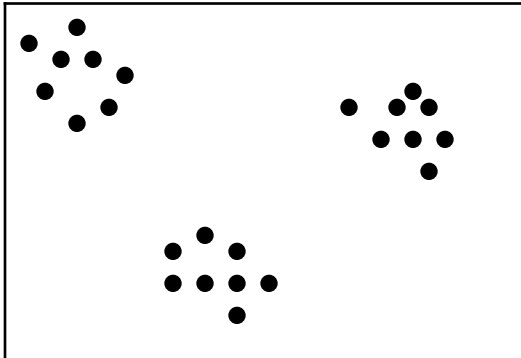
- Randomly reinitialize
  - ▣ take randomly selected  $k$  instances as the initial  $\mathbf{m}_i$
  - ▣ 1) calculate the mean of all data; 2) add small random vectors to the mean to get the  $k$  initial  $\mathbf{m}_i$ .

# Bad cases for k-means

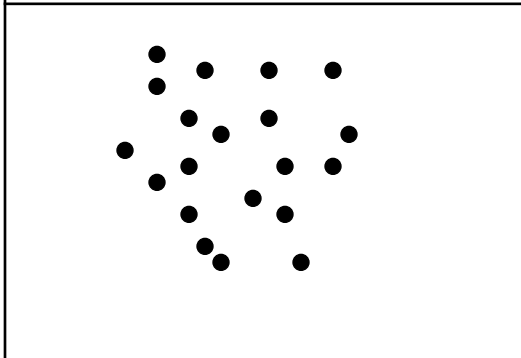
- Clusters may overlap
- Some clusters may be “wider” than others



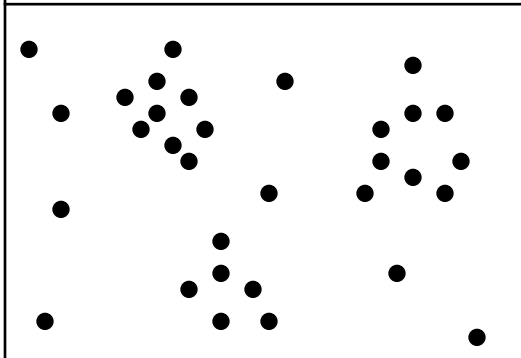
# Unsupervised Learning



Sometimes easy



Sometimes impossible



and sometimes in between

# K-means clustering using intensity or color

Image



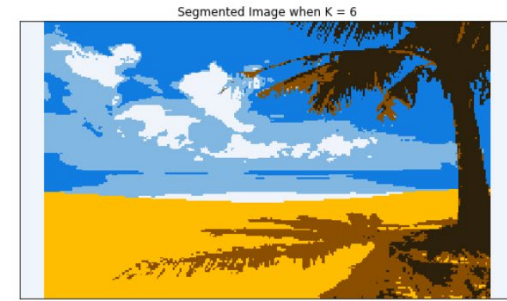
Clusters on intensity



Clusters on color



# Choosing $k$



- Defined by the application, e.g., image segmentation or color quantization
- Plot data (projection to low dimension) and check for clusters
- Add one at a time until small change in reconstruction error
- Cross-validation