

# Machine Learning and Data Analytics

## ME 5013- Fall 2019

### Lectures 16

- K-means clustering



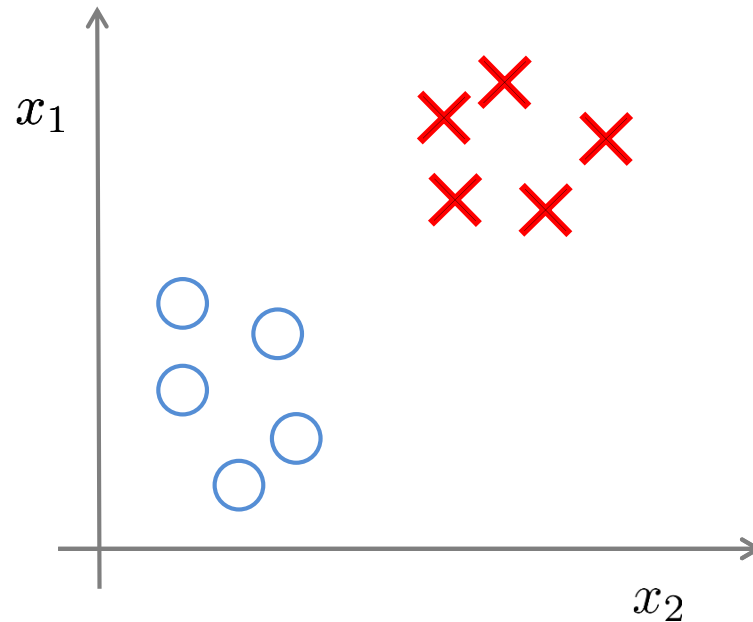
The University of Texas at San Antonio™

Adel Alaeddini, PhD

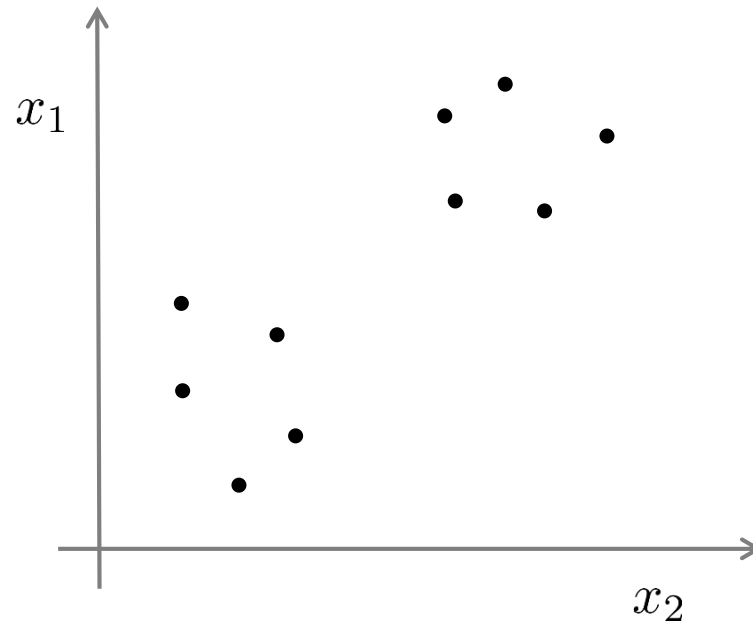
Associate Professor of Mechanical Engineering

Advanced Data Engineering Lab

[adel.alaeddini@utsa.edu](mailto:adel.alaeddini@utsa.edu)

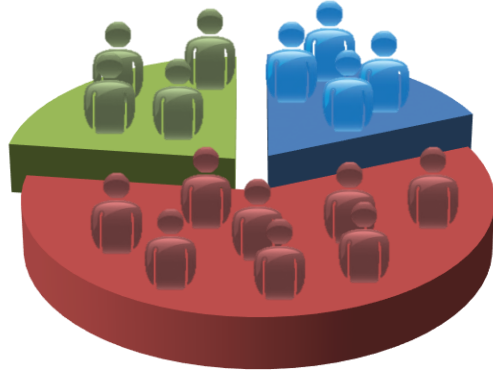


Training set:  $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), (x^{(3)}, y^{(3)}), \dots, (x^{(n)}, y^{(n)})\}$

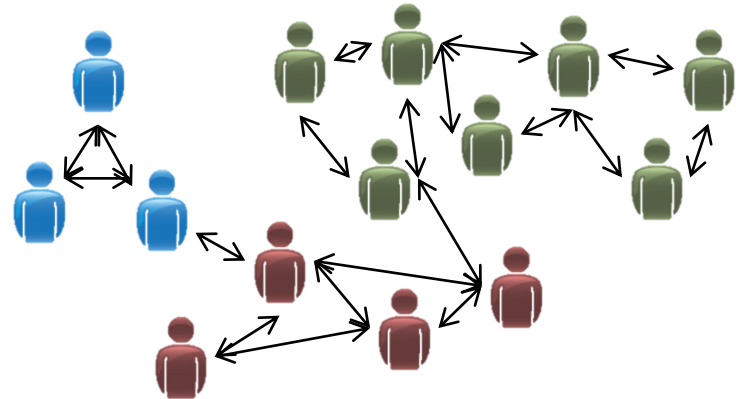


Training set:  $\{x^{(1)}, x^{(2)}, x^{(3)}, \dots, x^{(n)}\}$

## Applications of clustering



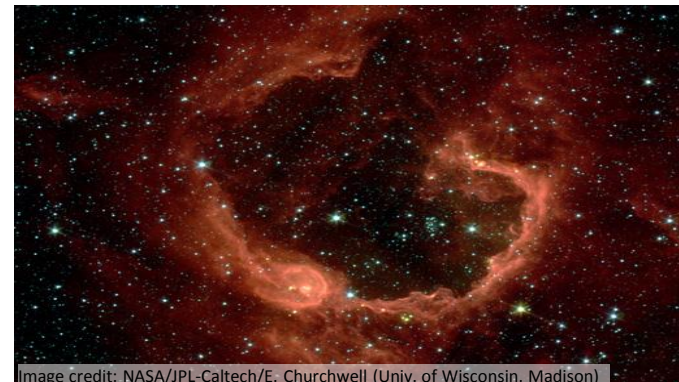
Market segmentation



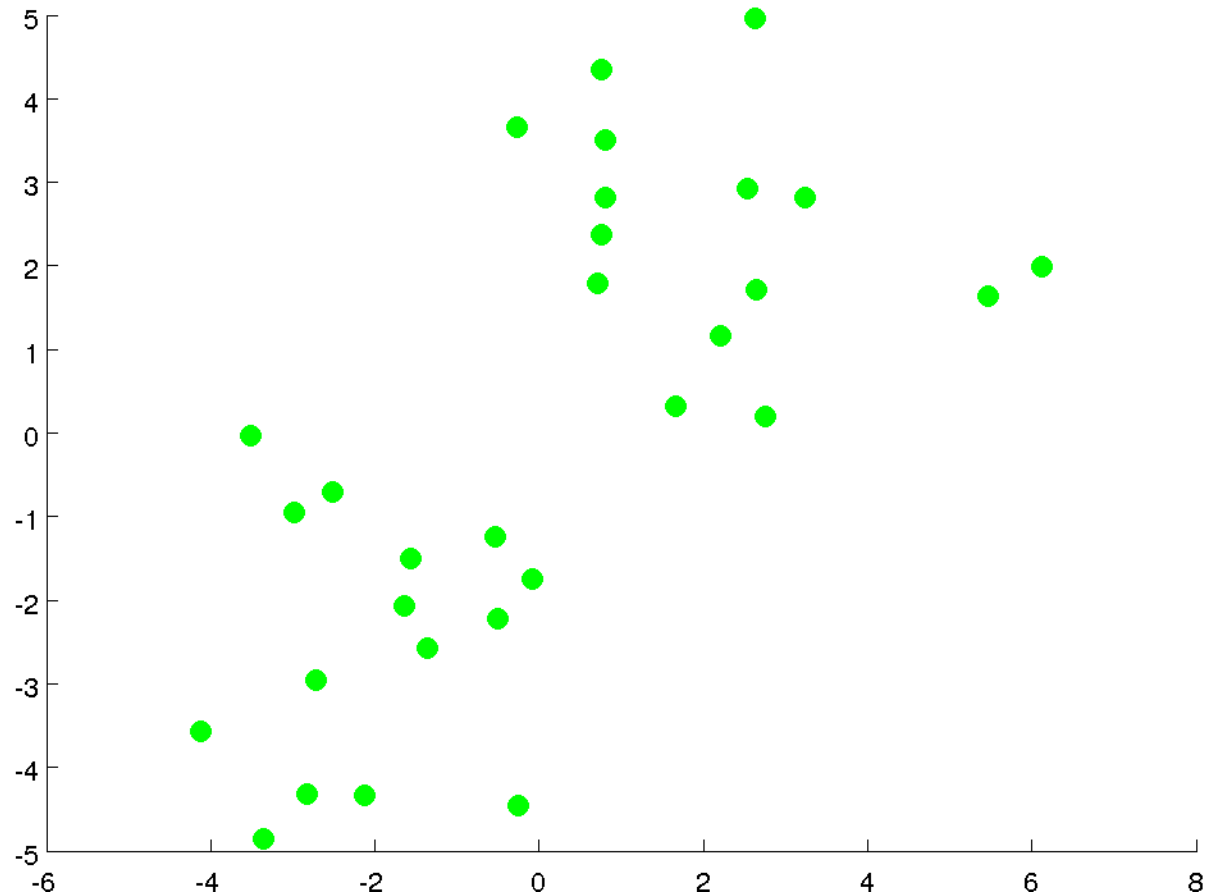
Social network analysis



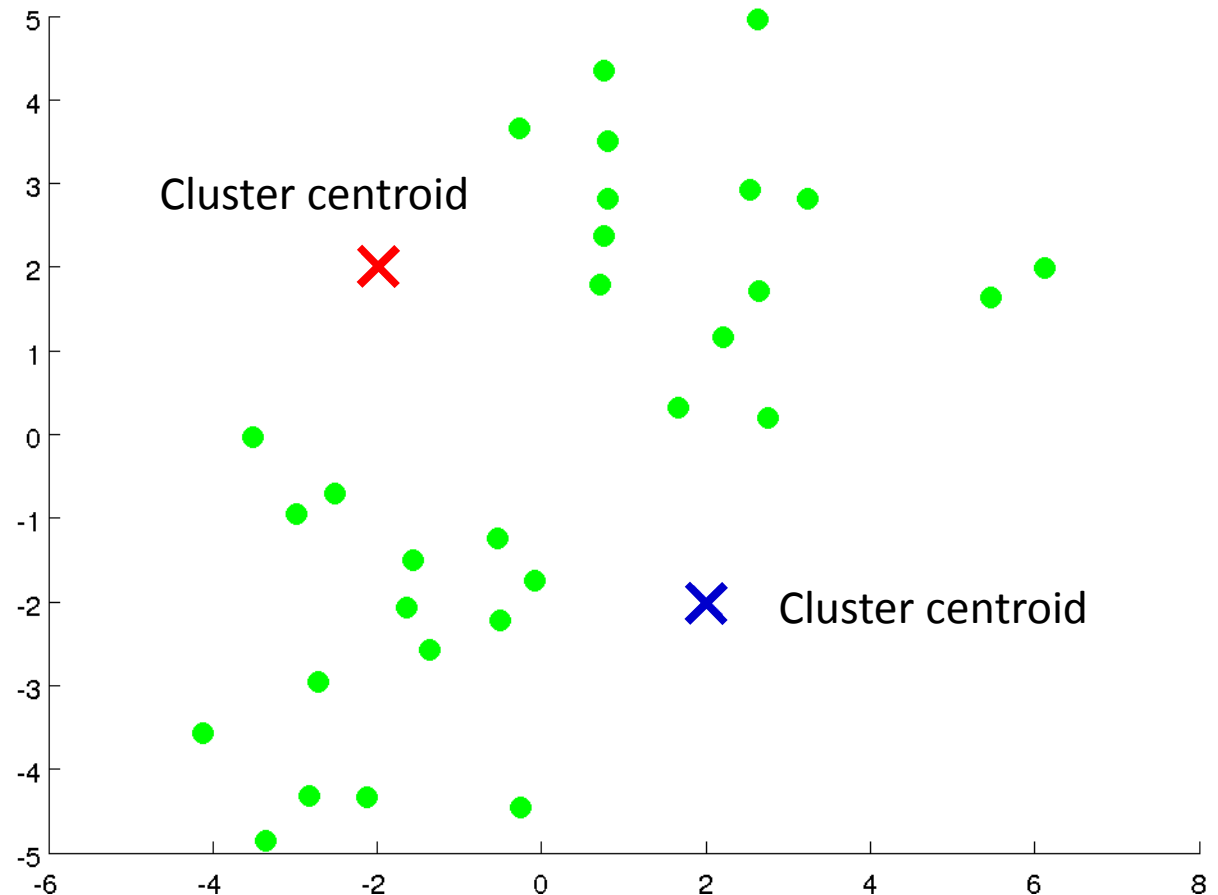
Organize computing clusters



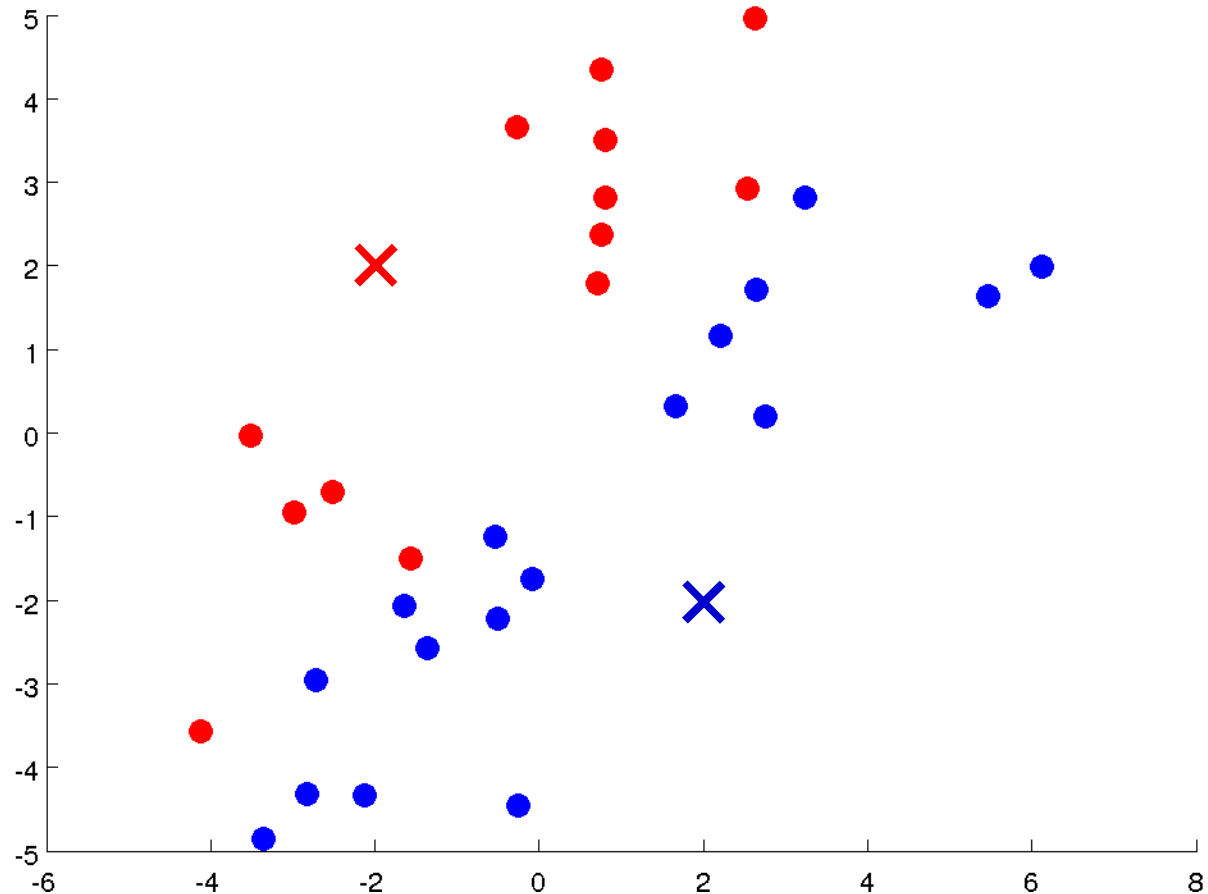
Astronomical data analysis



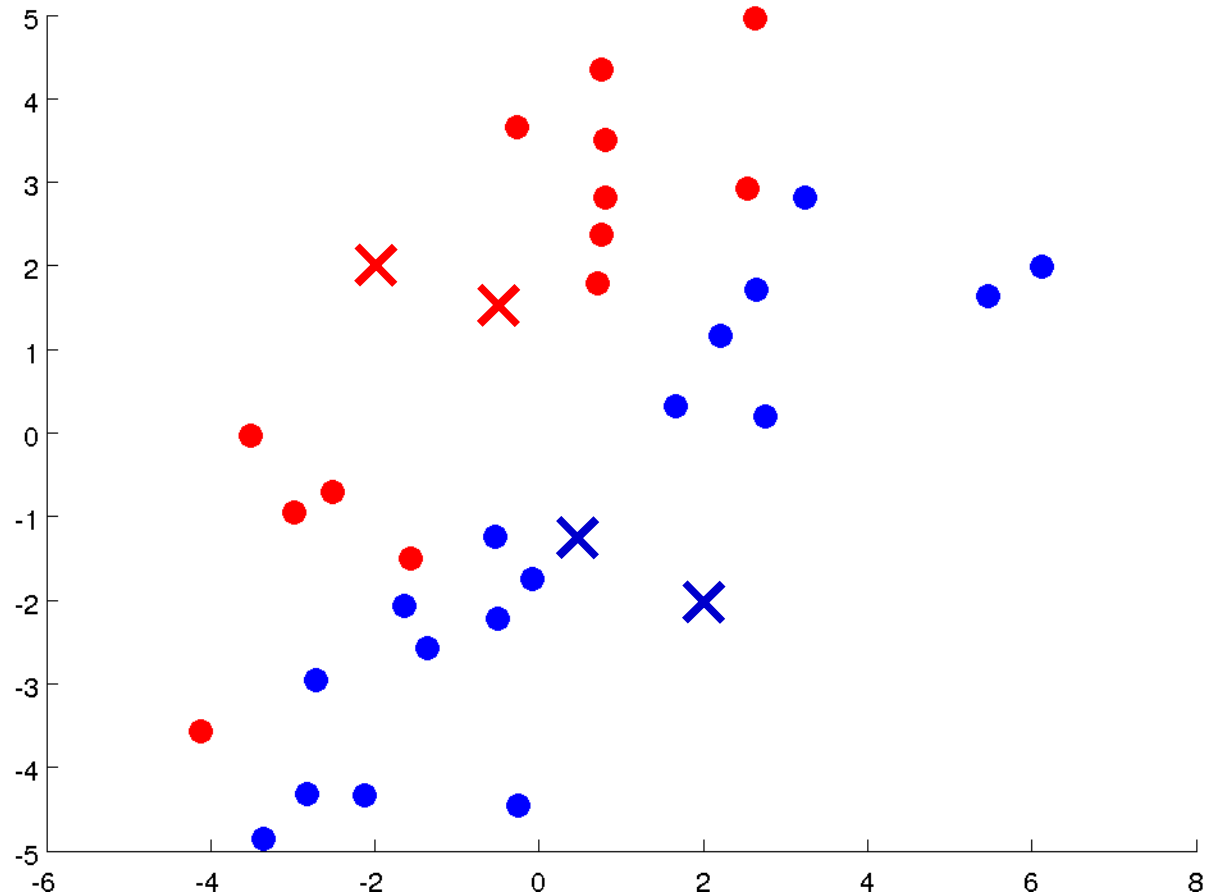
- Random selection of cluster centroids
- Loop
  - Cluster assignment
  - Move centroids



- Cluster assignment

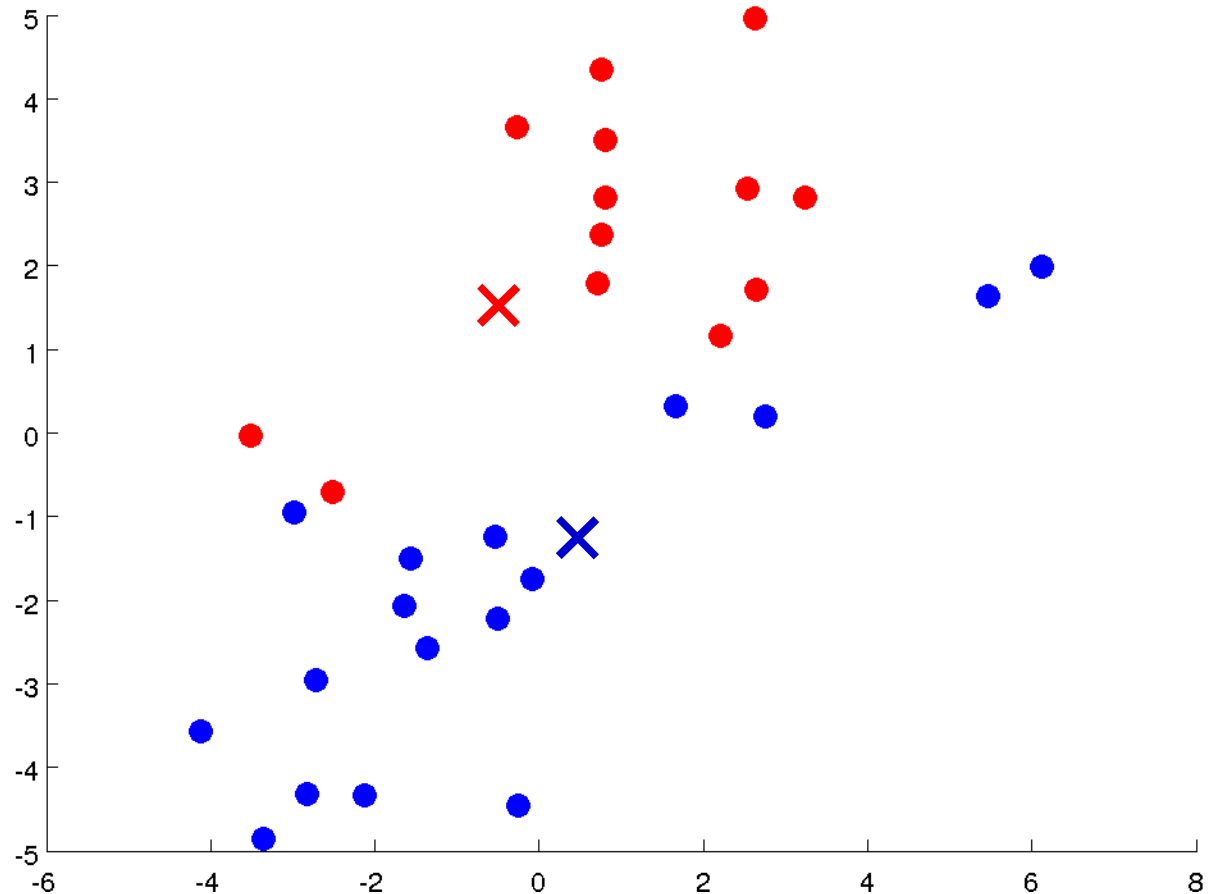


- Move centroids

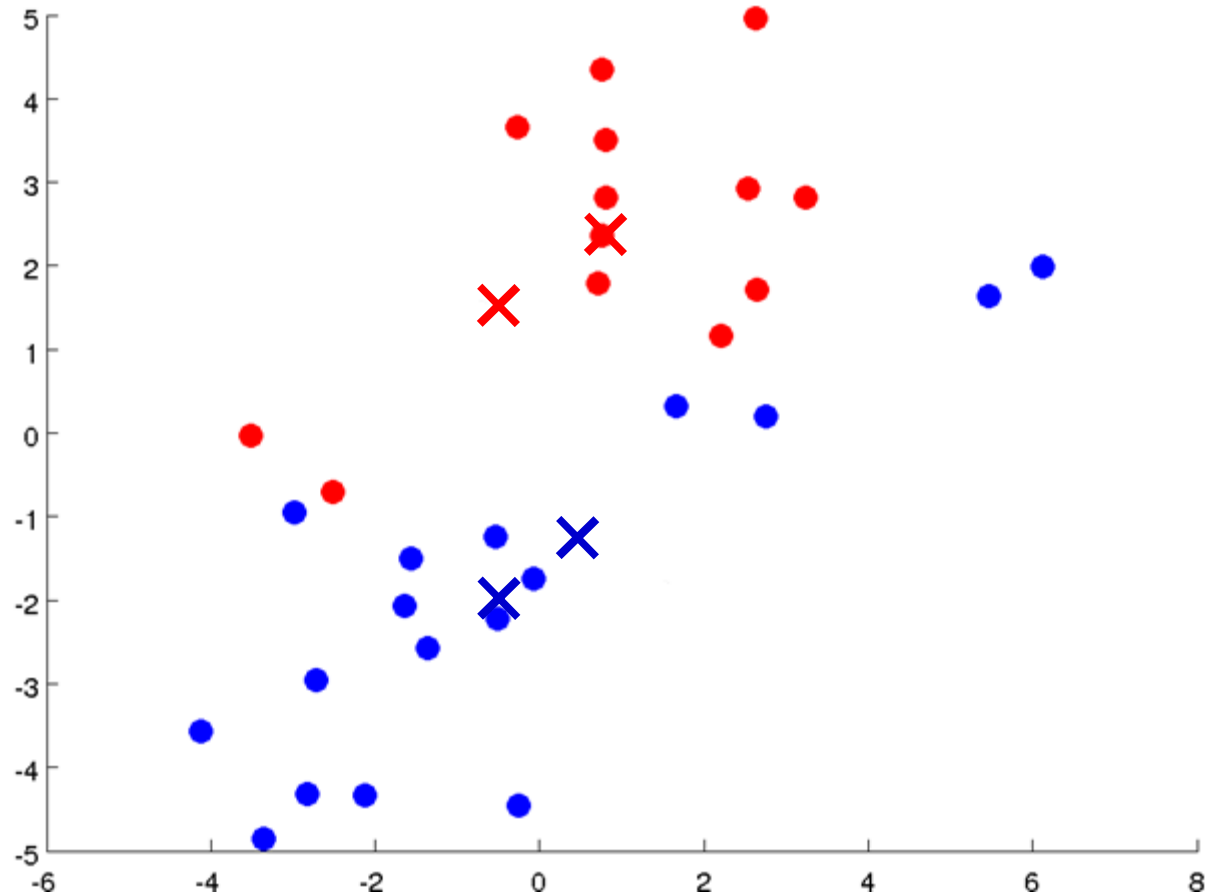




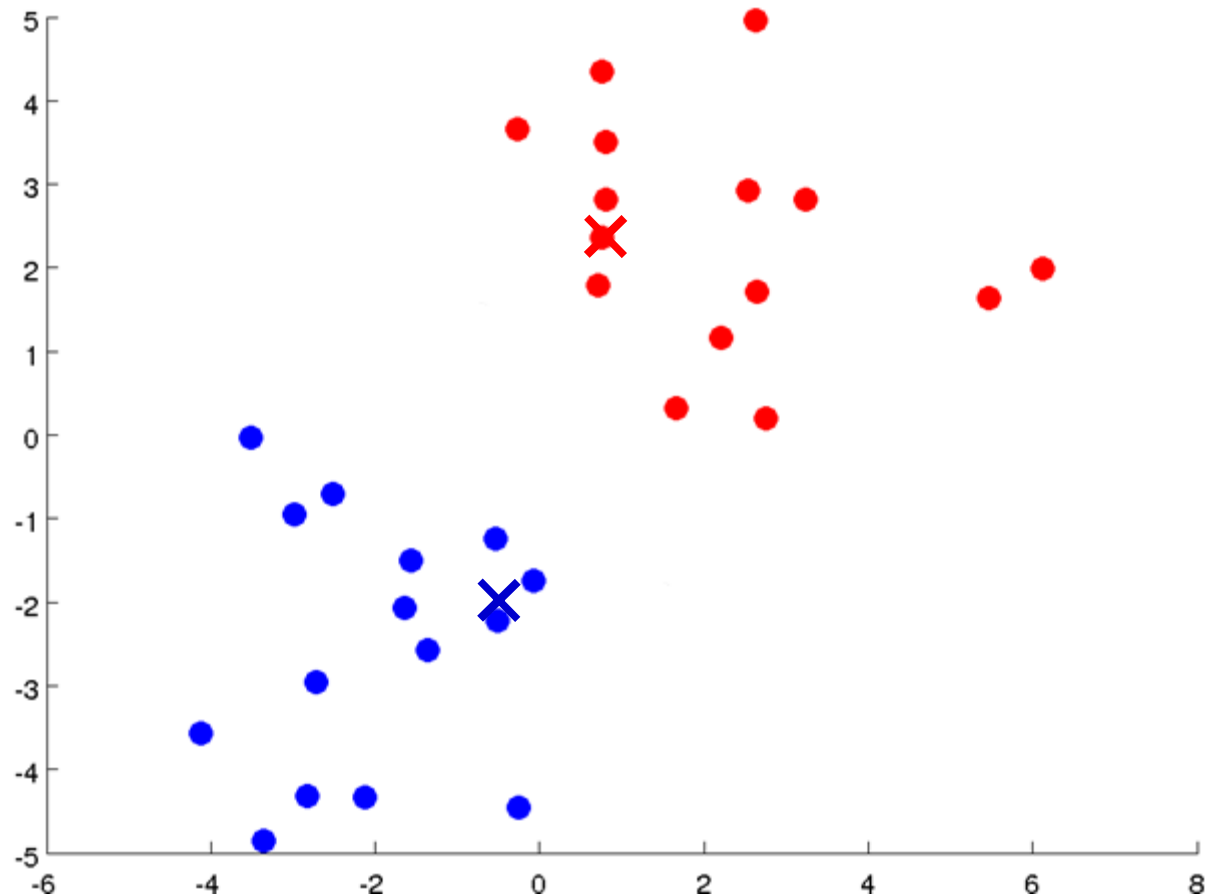
- Cluster assignment



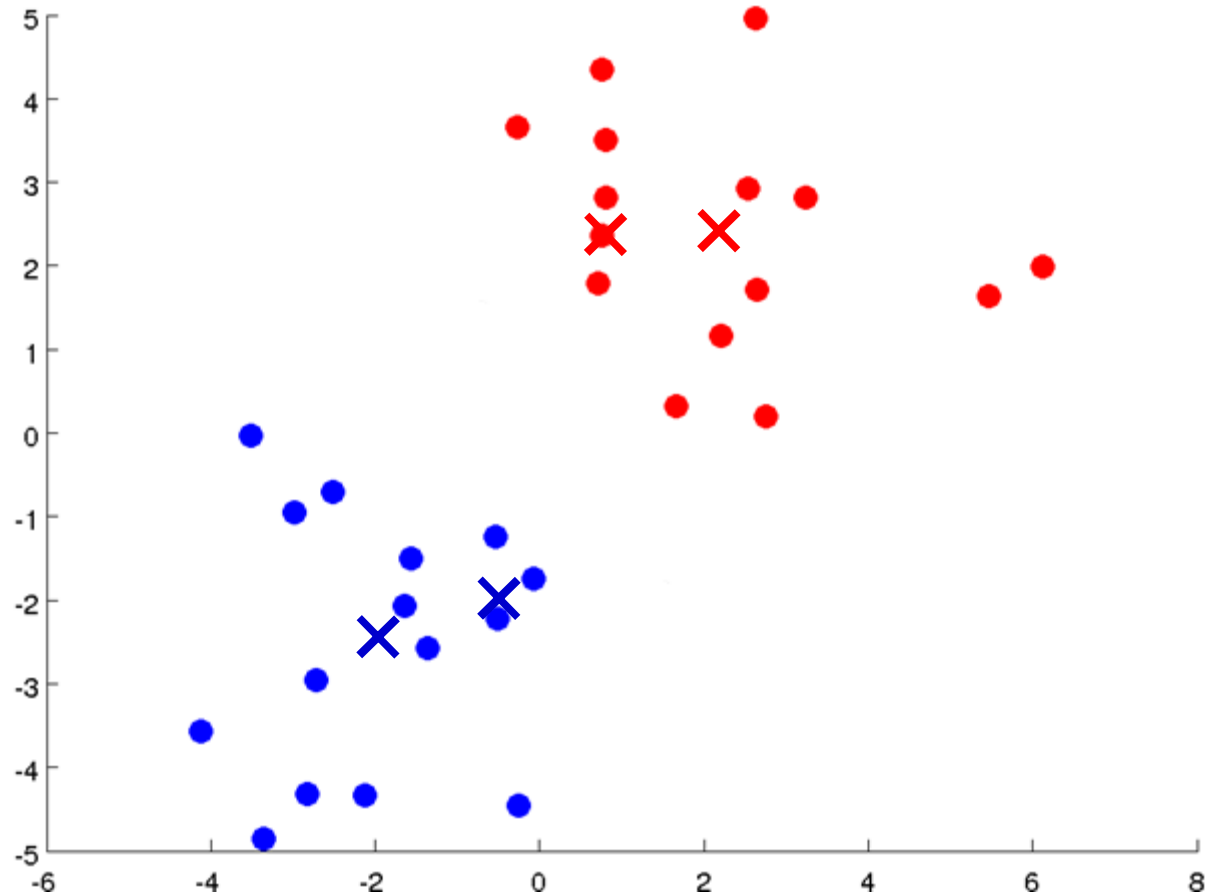
- Move centroids



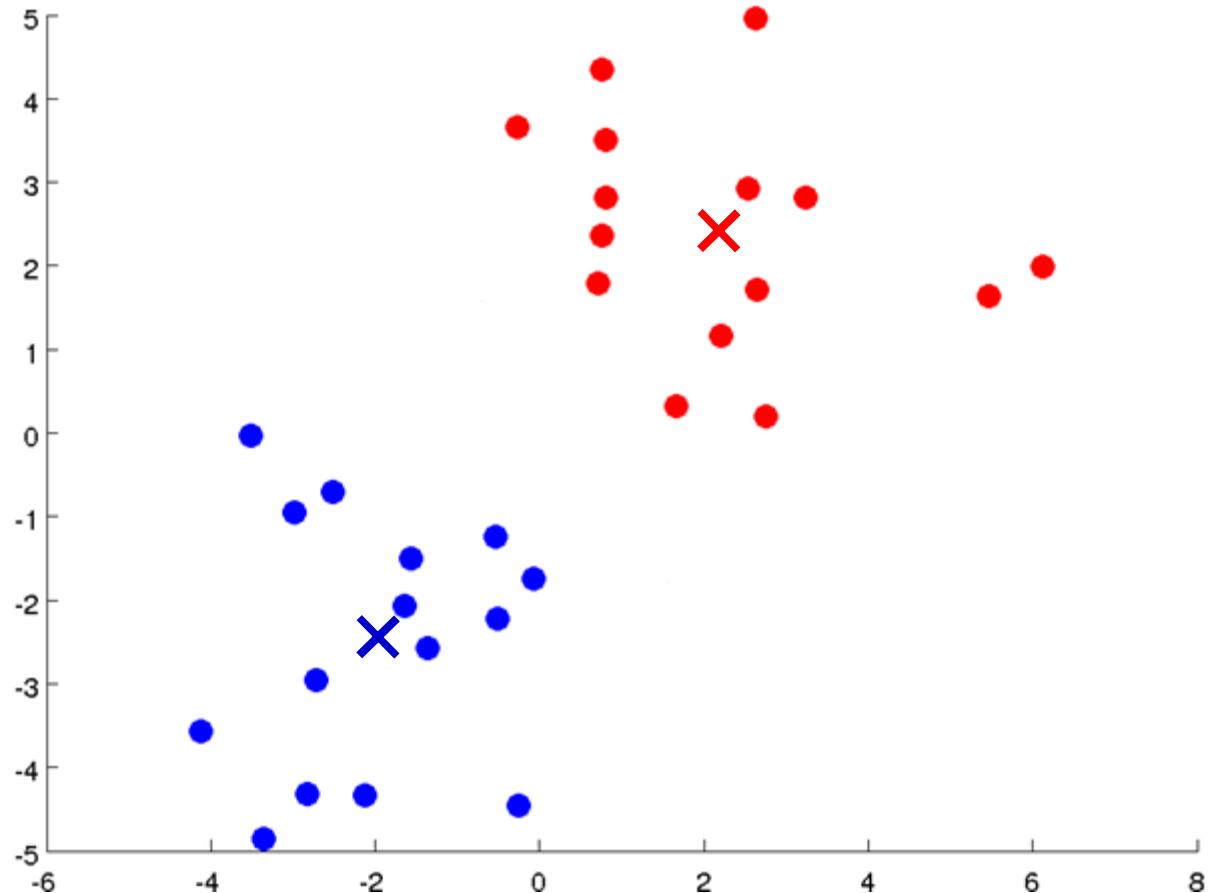
- Cluster assignment



- Move centroids



- Cluster assignment



## K-means algorithm

Input:

- $K$  (number of clusters)
- Training set  $\{x^{(1)}, x^{(2)}, \dots, x^{(n)}\}$

$x^{(i)} \in \mathbb{R}^p$  (drop  $x_0 = 1$  convention)

## K-means algorithm

Randomly initialize  $K$  cluster centroids  $\mu_1, \mu_2, \dots, \mu_K \in \mathbb{R}^p$

Repeat {

for  $i = 1$  to  $n$

$c^{(i)} :=$  index (from 1 to  $K$ ) of cluster centroid  
closest to  $x^{(i)}$

**Cluster  
assignment**

$$c^{(i)} = \operatorname{Argmin}_k \|x^{(i)} - \mu_k\|^2$$

for  $k = 1$  to  $K$

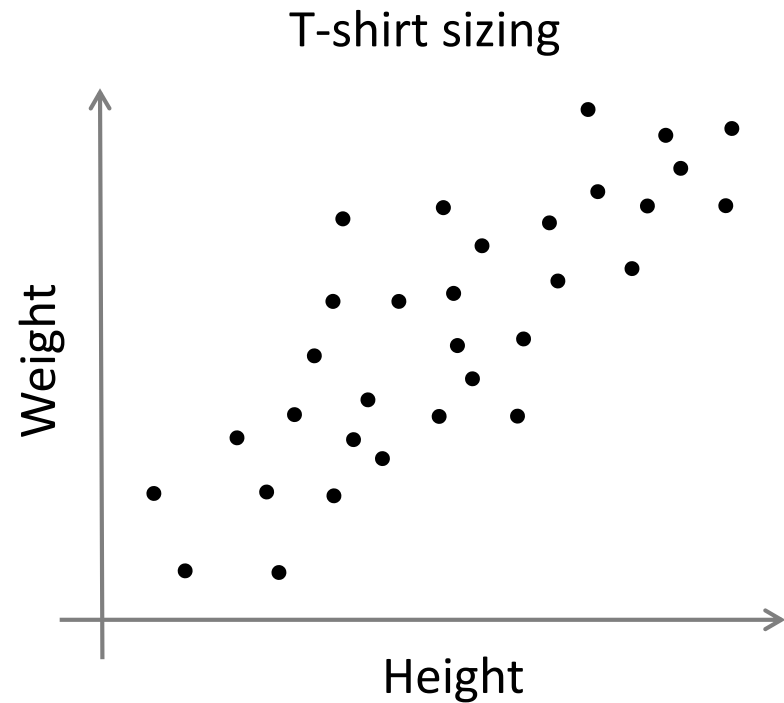
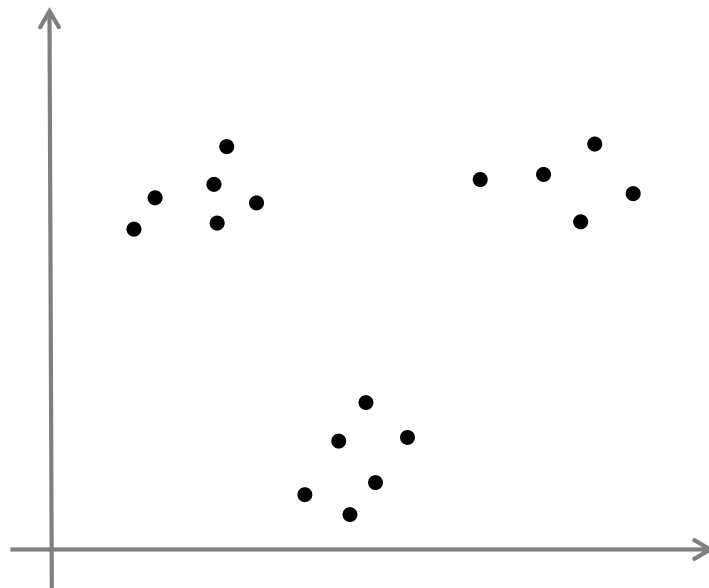
$\mu_k :=$  average (mean) of points assigned to cluster

**Move  
centroid**

$$\mu_k = \operatorname{Mean}_{c^{(i)}=k} x^{(i)}$$

}

## K-means for non-separated clusters





## K-means optimization objective

$c^{(i)}$  = index of cluster  $(1, 2, \dots, K)$  to which example  $x^{(i)}$  is currently assigned

$\mu_k$  = cluster centroid  $k$  ( $\mu_k \in \mathbb{R}^p$ )

$\mu_{c^{(i)}}$  = cluster centroid of cluster to which example  $x^{(i)}$  has been assigned

Optimization objective:

$$J(c^{(1)}, \dots, c^{(n)}, \mu_1, \dots, \mu_K) = \frac{1}{n} \sum_{i=1}^n \|x^{(i)} - \mu_{c^{(i)}}\|^2$$

$$\min_{\substack{c^{(1)}, \dots, c^{(n)}, \\ \mu_1, \dots, \mu_K}} J(c^{(1)}, \dots, c^{(n)}, \mu_1, \dots, \mu_K)$$

## K-means algorithm

Randomly initialize  $K$  cluster centroids  $\mu_1, \mu_2, \dots, \mu_K \in \mathbb{R}^p$

Repeat {

for  $i = 1$  to  $n$

$c^{(i)} := \text{index (from 1 to } K \text{) of cluster centroid}$   
 $\text{closest to } x^{(i)}$

**Cluster  
assignment**

Minimize  $J(.)$  w.r.t.  $c^{(i)}$   
 (holding  $\mu_k$  fixed)

for  $k = 1$  to  $K$

$\mu_k := \text{average (mean) of points assigned to cluster}$

**Move  
centroid**

$\mu_k = \text{Mean}_{c^{(i)}=k} x^{(i)}$

}

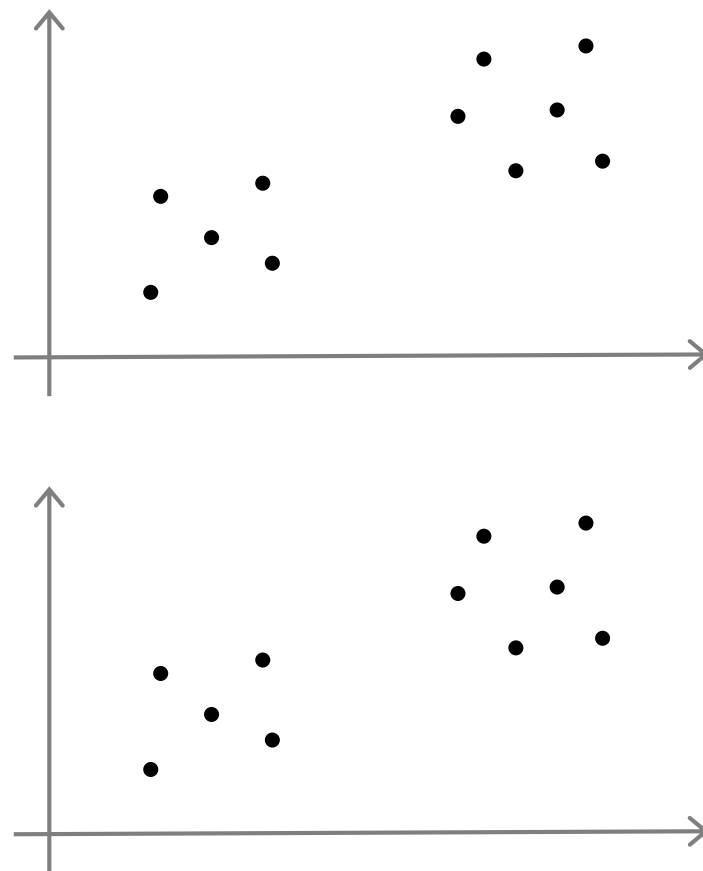
Minimize  $J(.)$   
 w.r.t.  $\mu_k$  (holding  $c^{(i)}$  fixed)

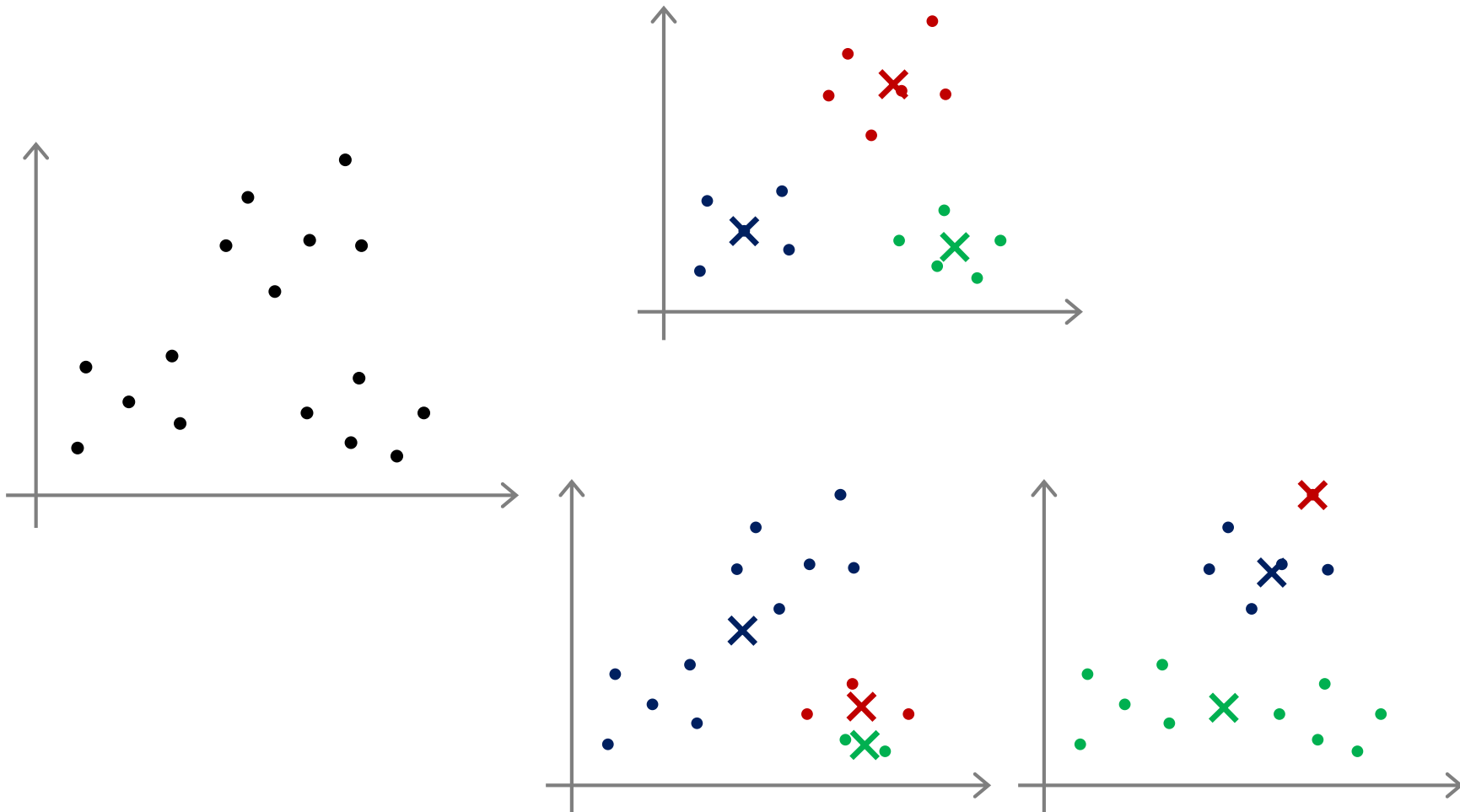
## Random initialization

Should have  $K < n$

Randomly pick  $K$  training examples.

Set  $\mu_1, \dots, \mu_K$  equal to these  $K$  examples.





For  $i = 1$  to 100 {

Randomly initialize K-means.

Run K-means. Get  $c^{(1)}, \dots, c^{(n)}, \mu_1, \dots, \mu_K$

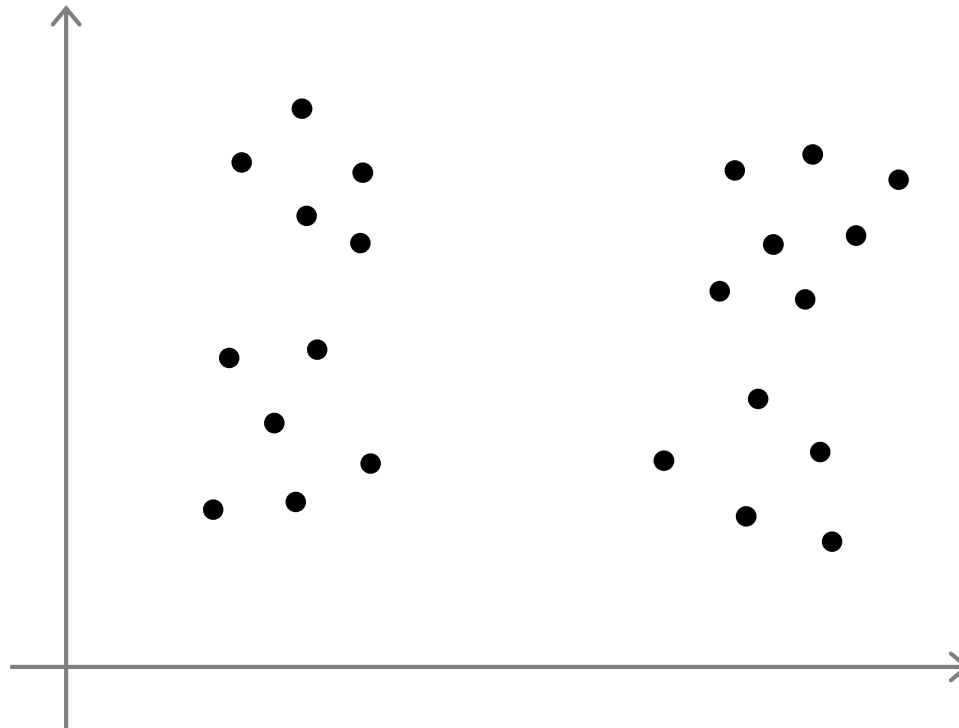
Compute cost function (distortion)

$$J(c^{(1)}, \dots, c^{(n)}, \mu_1, \dots, \mu_K)$$

}

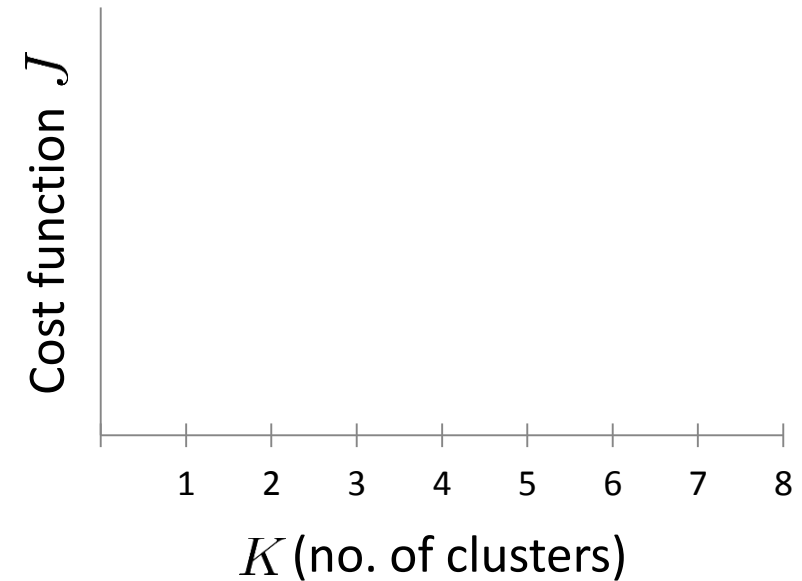
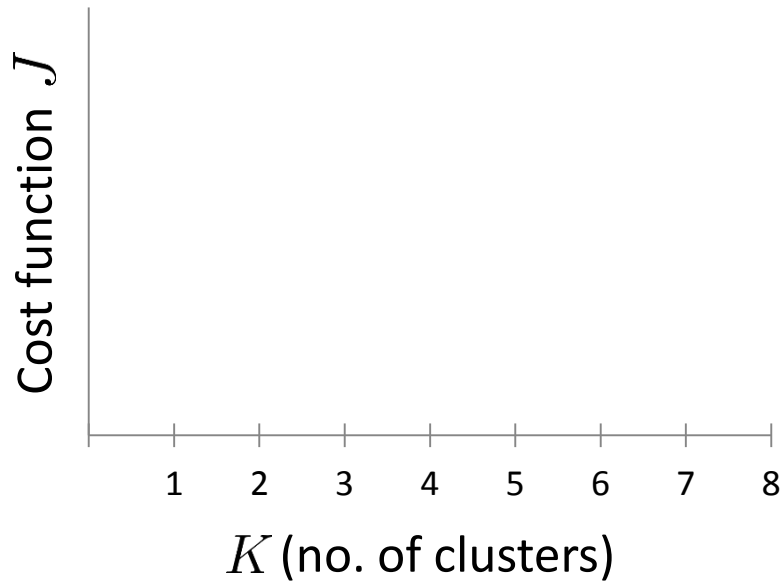
Pick clustering that gave lowest cost  $J(c^{(1)}, \dots, c^{(n)}, \mu_1, \dots, \mu_K)$

What is the right value of K?



## Choosing the value of $K$

Elbow method:



Sometimes, you're running K-means to get clusters to use for some later/downstream purpose. Evaluate K-means based on a metric for how well it performs for that later purpose.

E.g.

