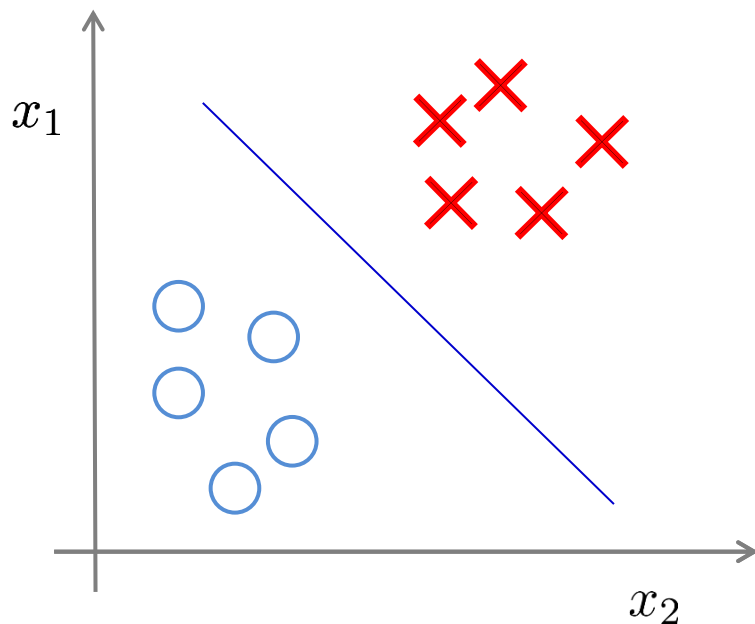# Clustering

Machine Learning (AIM 5002-41)

Joon Hee Choi
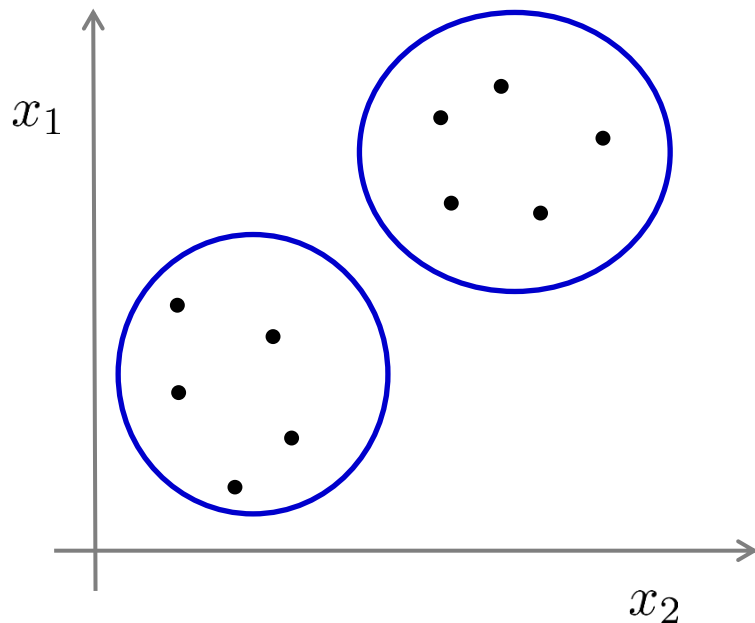Sungkyunkwan University

# Clustering:
# Unsupervised Learning

# Supervised learning



Training set: $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), (x^{(3)}, y^{(3)}), \ldots, (x^{(m)}, y^{(m)})\}$
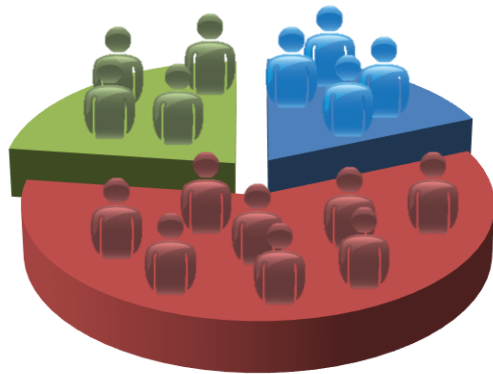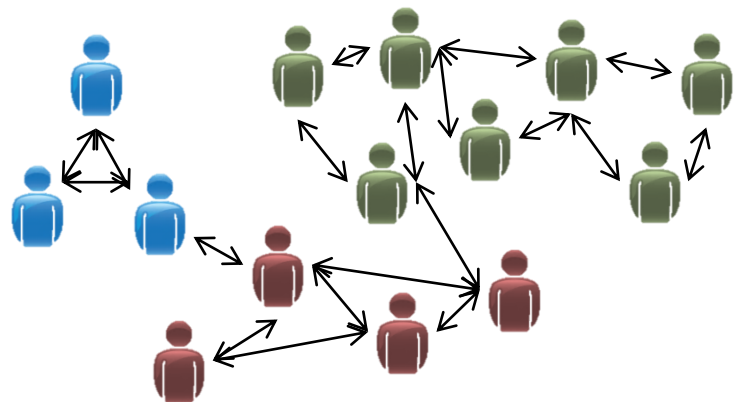
# Unsupervised learning



Clustering algorithm

Training set: $\{x^{(1)}, x^{(2)}, x^{(3)}, \dots, x^{(m)}\}$
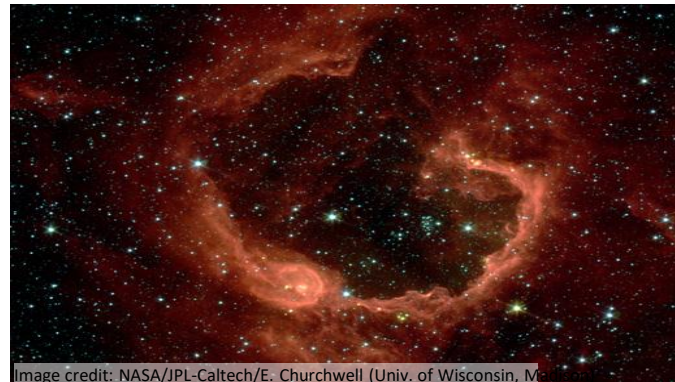
# Applications of clustering
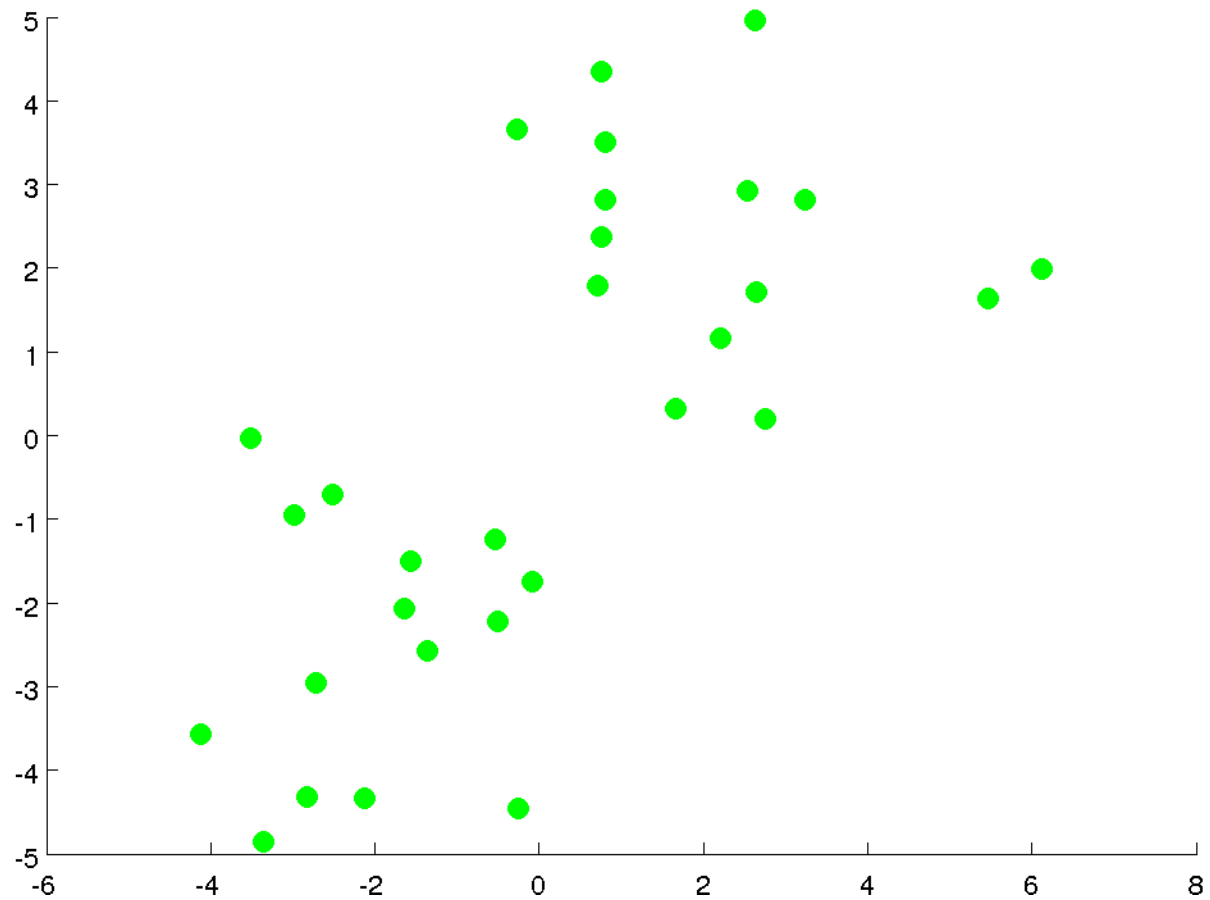


Market segmentation

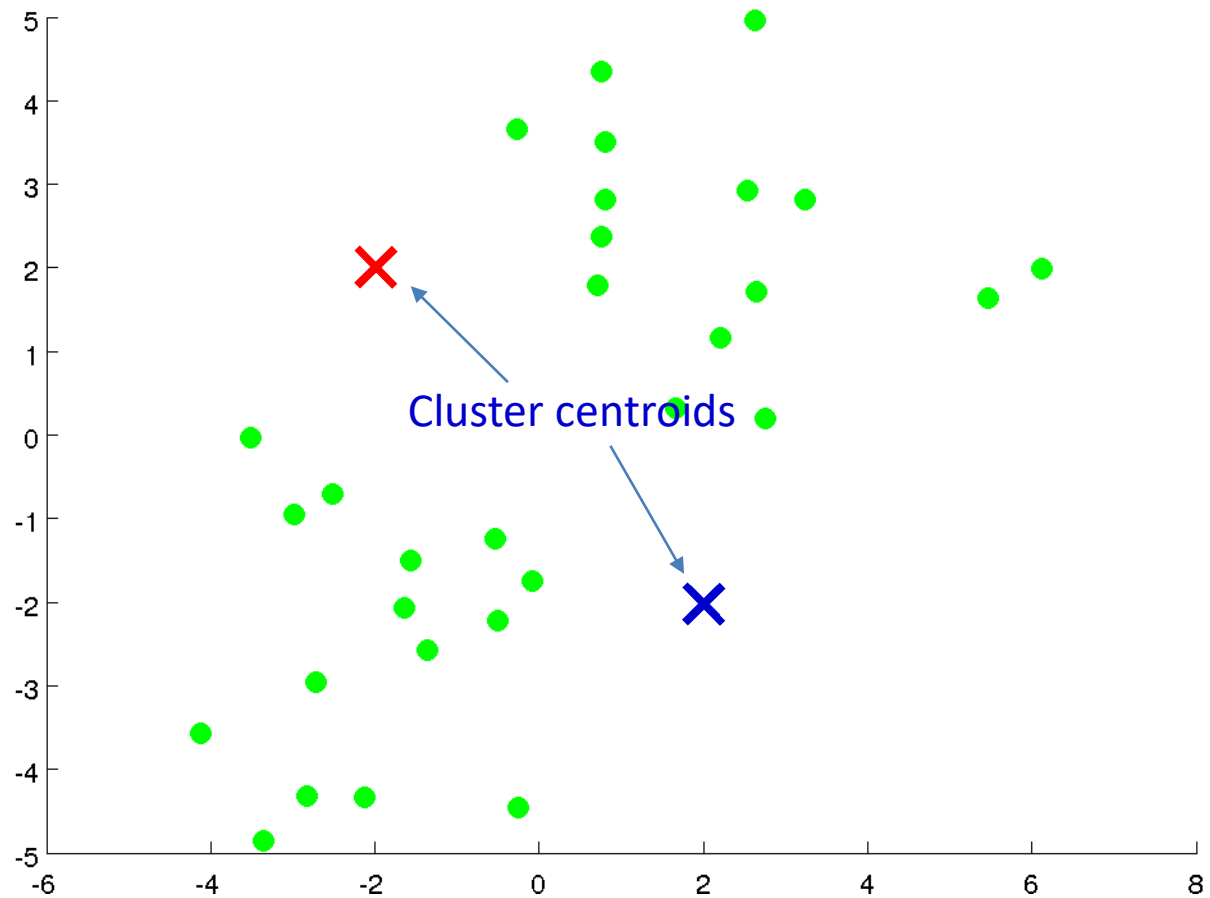Social network analysis

Organize computing clusters

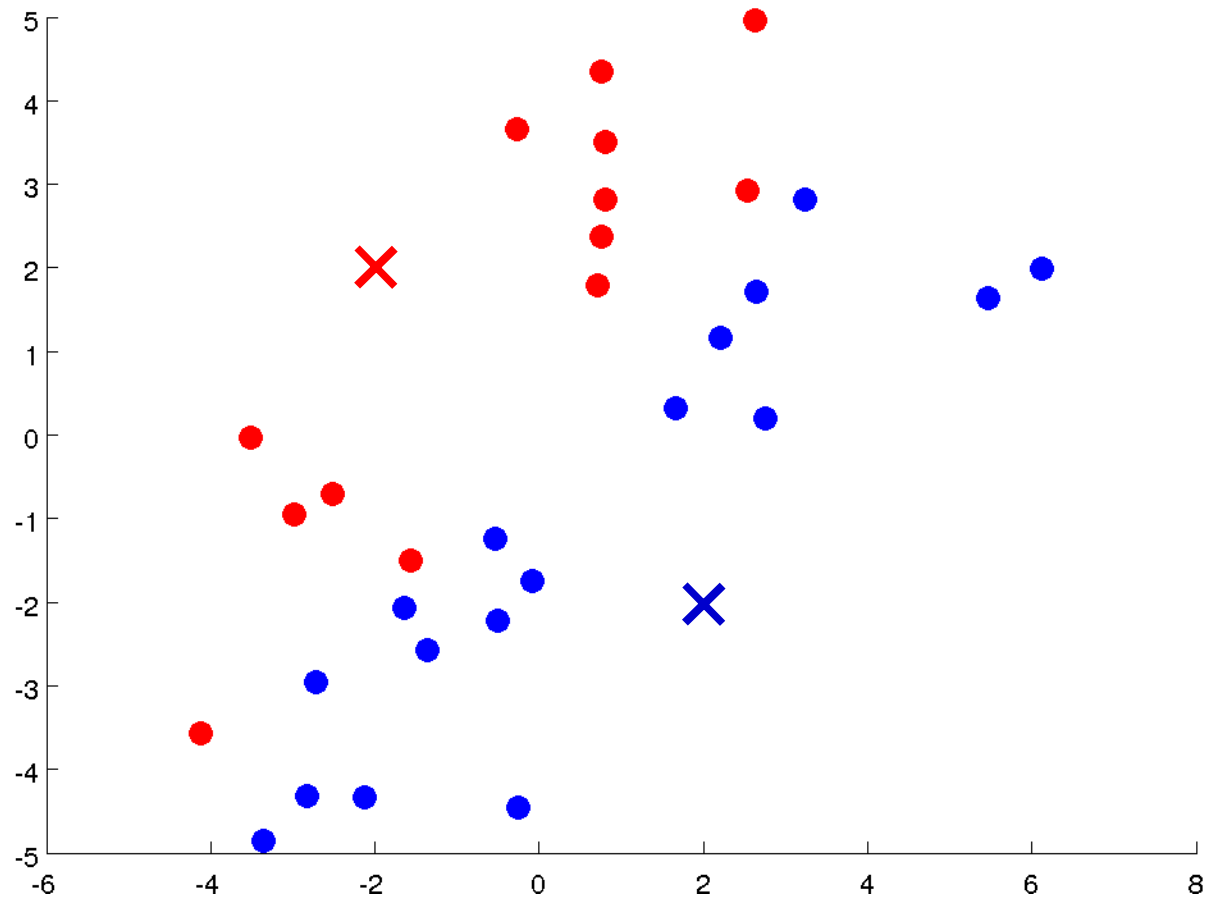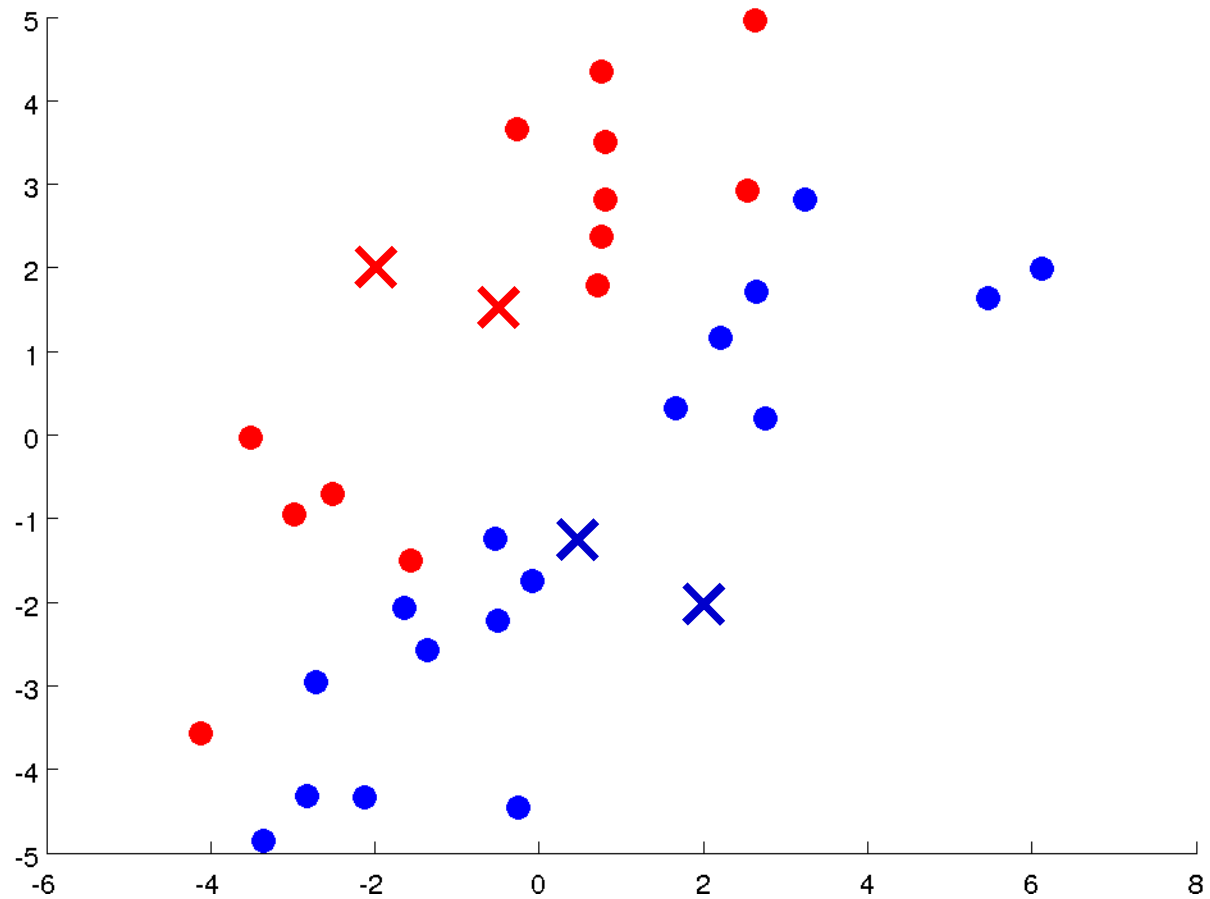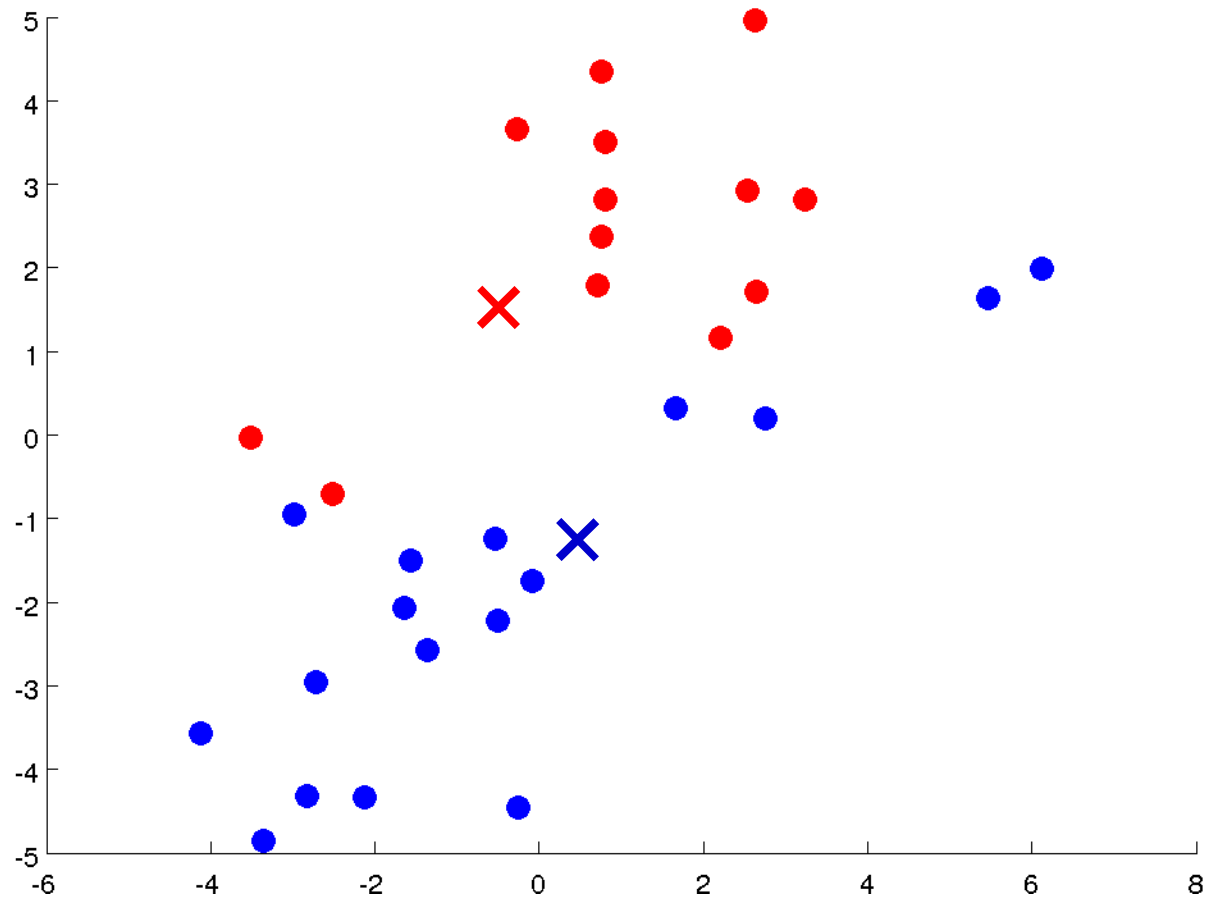Astronomical data analysis

Image credit: NASA/JPL-Caltech/E. Churchwell (Univ. of Wisconsin, M...

Andrew Ng

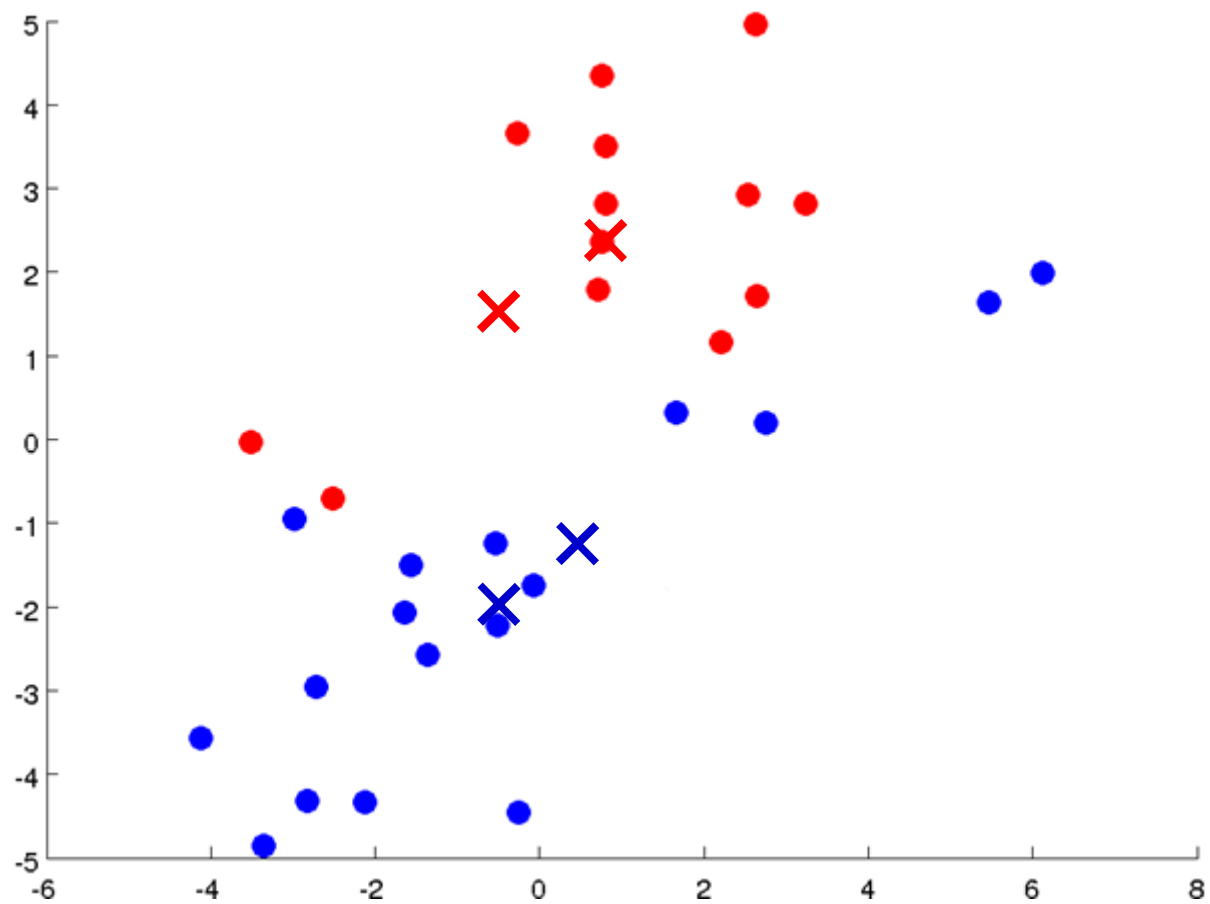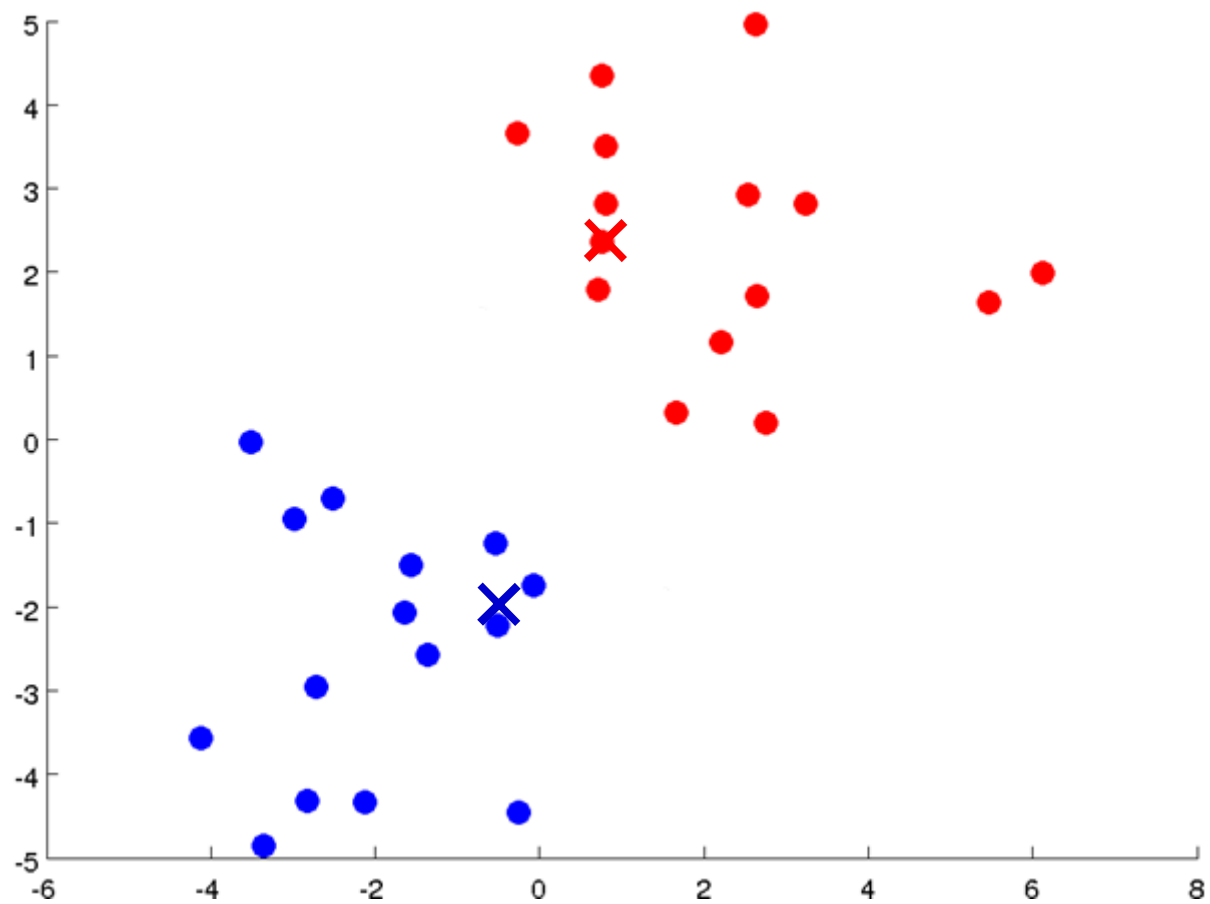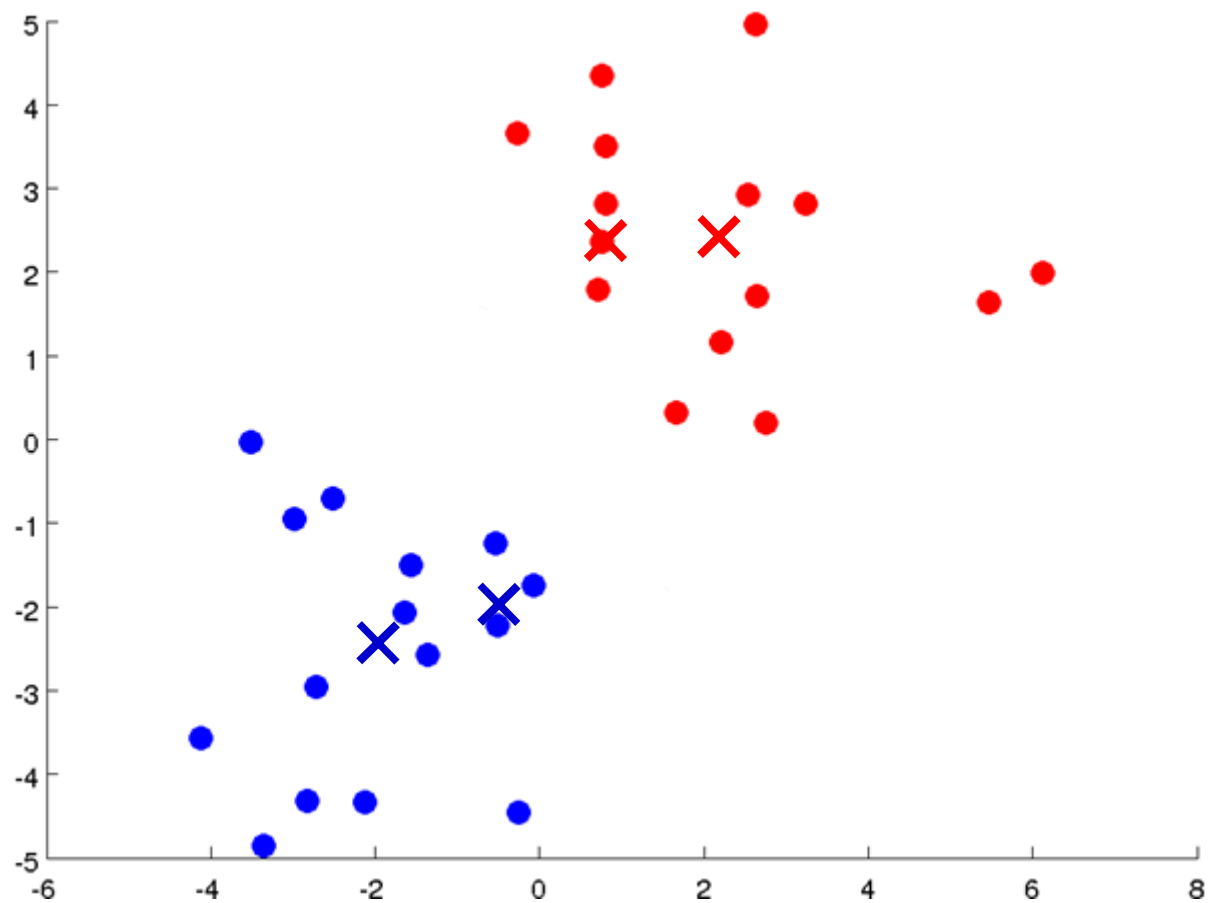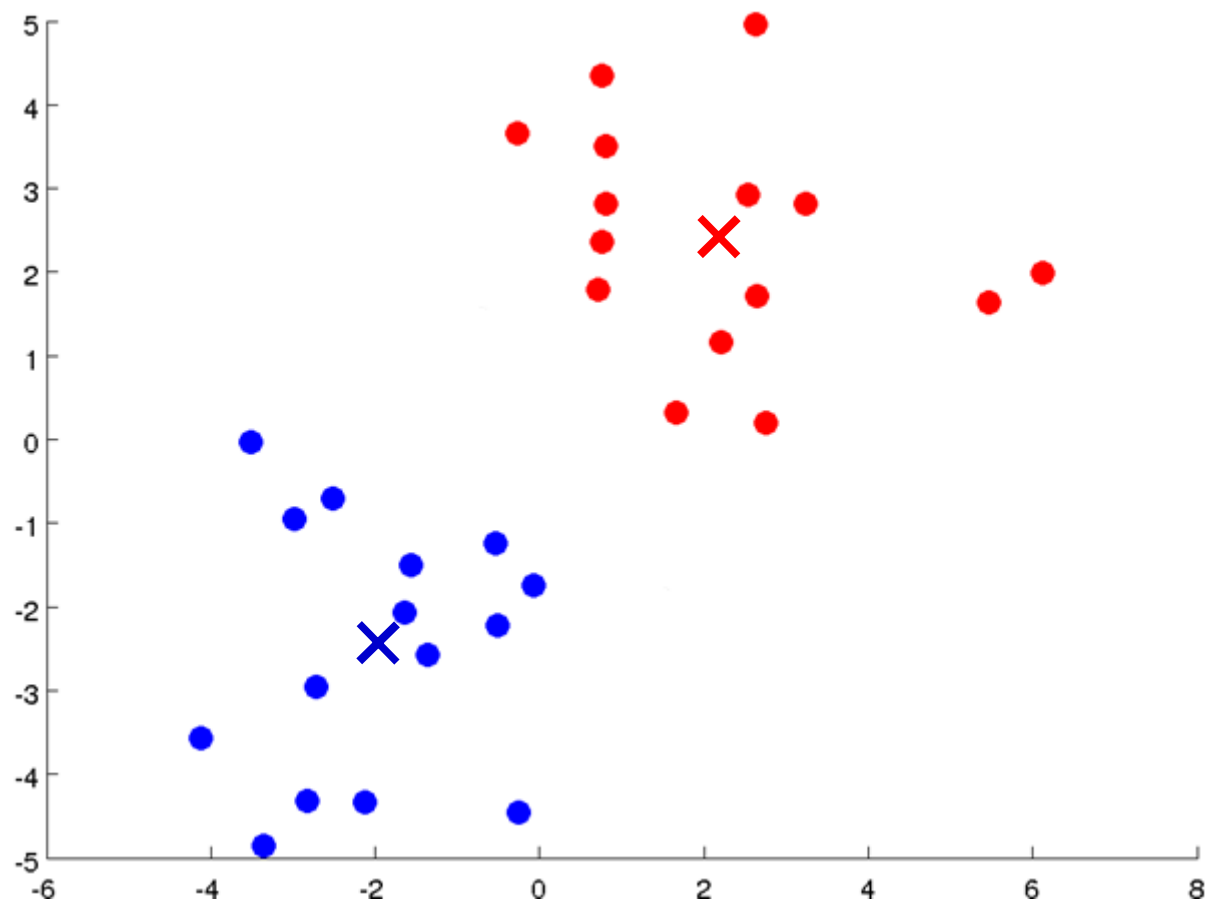# Clustering:
# K-means algorithm

Cluster centroids

# K-means algorithm

Input:

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

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

# K-means algorithm

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

Repeat {

<span style="color:blue">Cluster assignment step</span>

    for $i$ = 1 to $m$

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

          closest to $x^{(i)}$      $c^{(i)} = \min_k \left\| x^{(i)} - \mu_k \right\|^2$

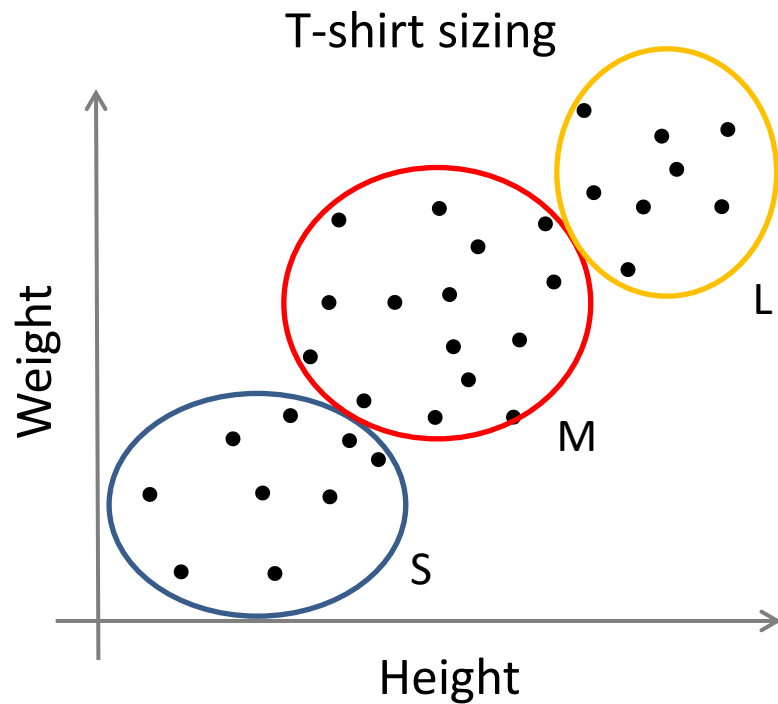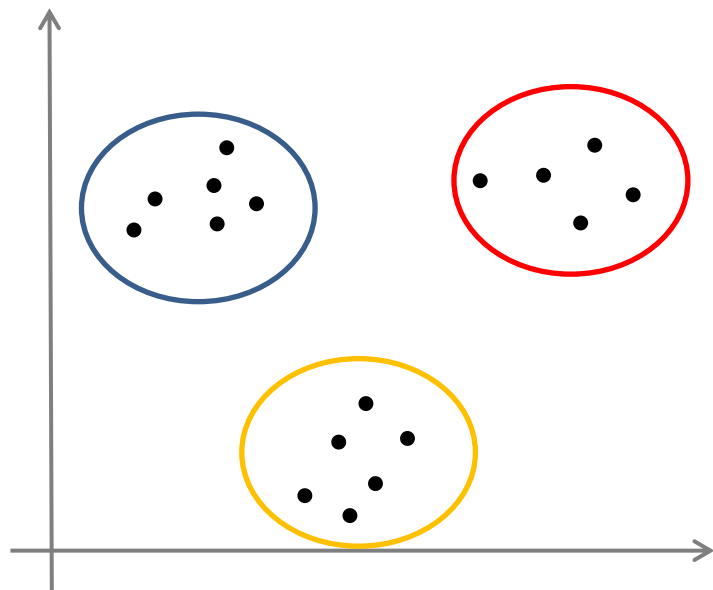<span style="color:blue">Move centroid</span>

    for $k$ = 1 to $K$

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

$$\mu_2 = \frac{1}{4}\left[x^{(1)} + x^{(5)} + x^{(6)} + x^{(10)}\right] \in \mathbb{R}^n$$

   }

# K-means for non-separated clusters



T-shirt sizing

Weight

Height

L

M

S

Andrew Ng

# Clustering:
# Optimization objective

Andrew Ng

# K-means optimization objective

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

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

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

Optimization objective:

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

$$\min_{\substack{c^{(1)}, \ldots, c^{(m)}, \\ \mu_1, \ldots, \mu_K}} J(c^{(1)}, \ldots, c^{(m)}, \mu_1, \ldots, \mu_K)$$

# K-means algorithm

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

Repeat {

min $J(\ldots)$ w.t $c^{(1)}, c^{(2)}, \ldots, c^{(m)}$
(holding $\mu_1, \mu_2, \ldots, \mu_k$ fixed)

Cluster assignment step

    for $i$ = 1 to $m$

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

           closest to $x^{(i)}$

Move centroid

    for $k$ = 1 to $K$

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

}

min $J(\ldots)$ w.t $\mu_1, \mu_2, \ldots, \mu_k$

Andrew Ng

# Clustering:
# Random initialization

**K-means algorithm**

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

Repeat {
     for $i$ = 1 to $m$
        $c^{(i)}$ := index (from 1 to $K$ ) of cluster centroid
            closest to $x^{(i)}$
     for $k$ = 1 to $K$
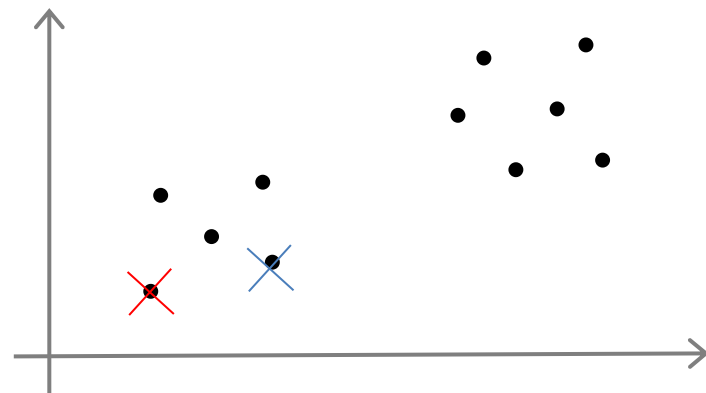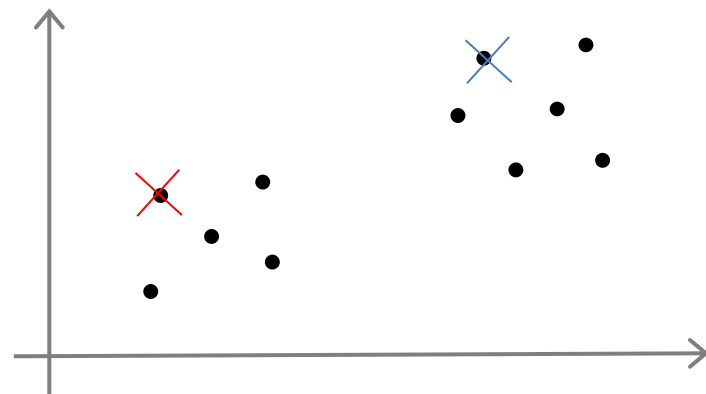        $\mu_k$  := average (mean) of points assigned to cluster $k$
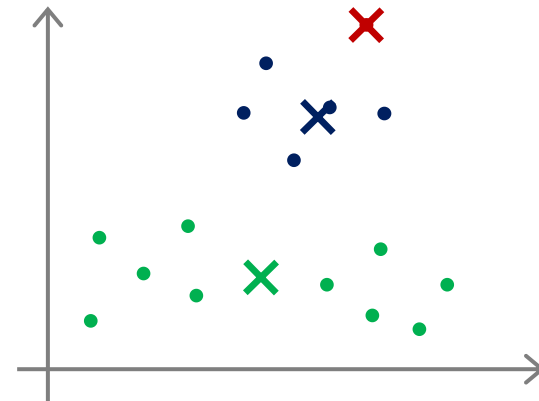}
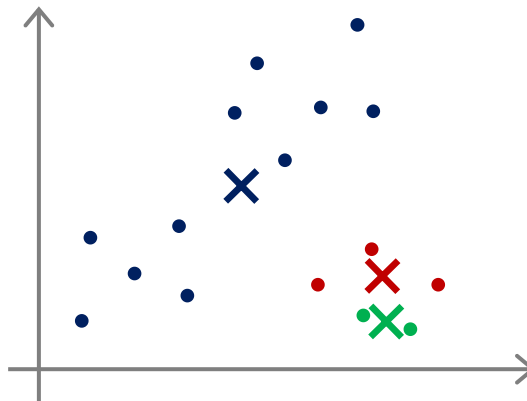
# Random initialization
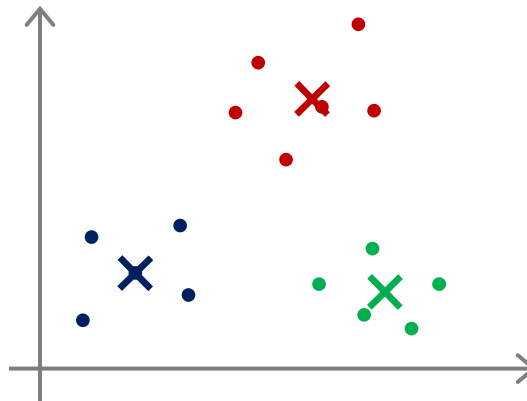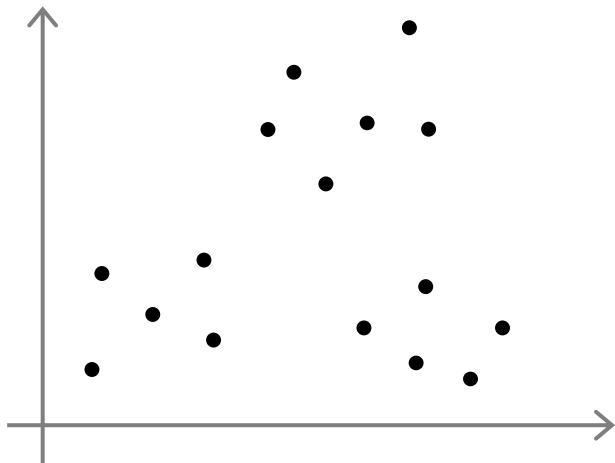
Should have $K < m$

Randomly pick $K$ training examples.

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

K=2

# Local optima

**Random initialization**

For i = 1 to 100 {

      Randomly initialize K-means.
      Run K-means. Get $c^{(1)}, \ldots, c^{(m)}, \mu_1, \ldots, \mu_K$.
      Compute cost function (distortion)
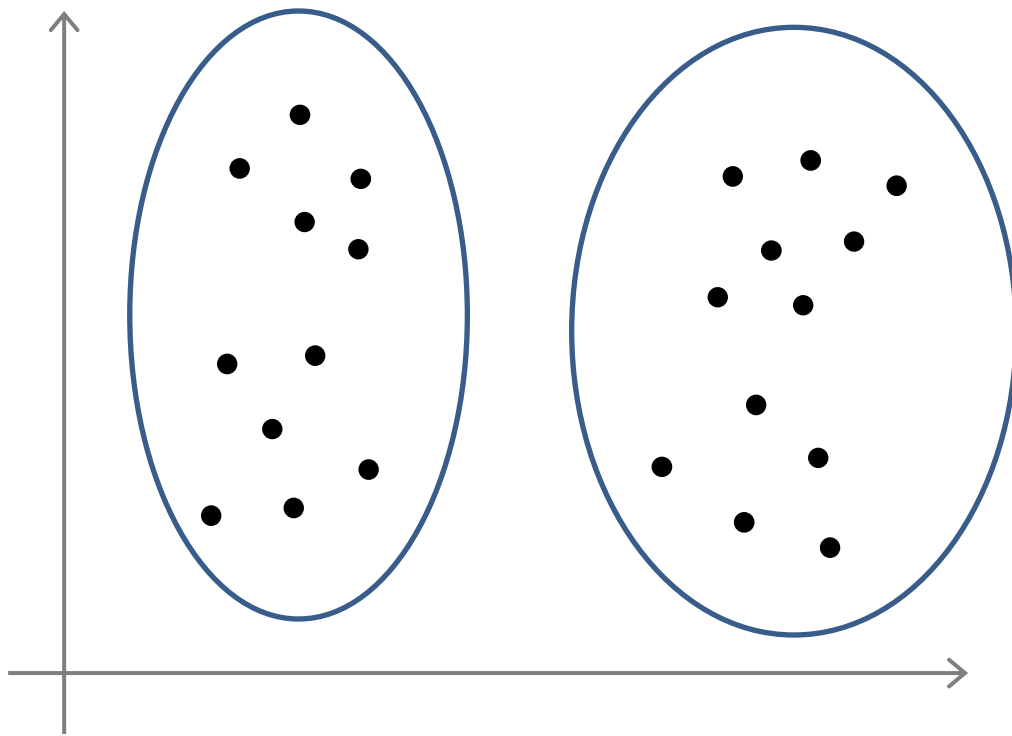        $J(c^{(1)}, \ldots, c^{(m)}, \mu_1, \ldots, \mu_K)$
      }

Pick clustering that gave lowest cost $J(c^{(1)}, \ldots, c^{(m)}, \mu_1, \ldots, \mu_K)$
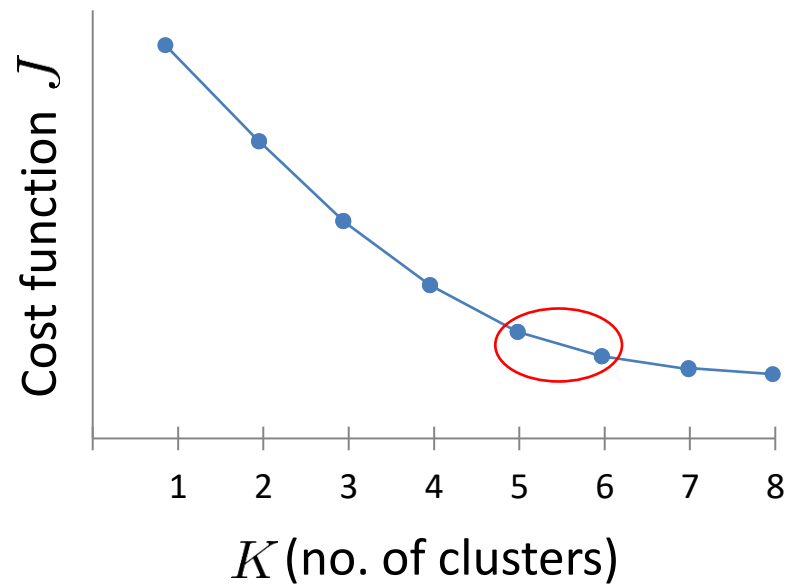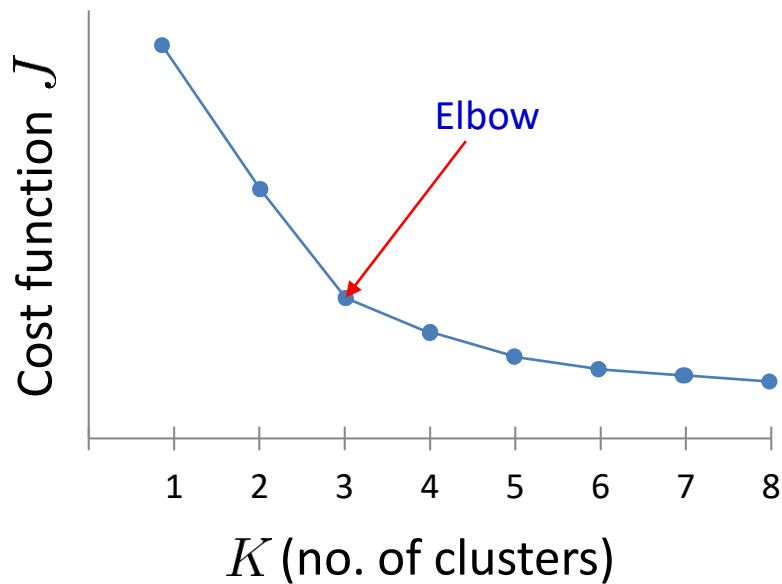
Clustering:
Choosing the number of clusters

Andrew Ng

# What is the right value of K?

# Choosing the value of K

Elbow method:

# Choosing the value of K

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.