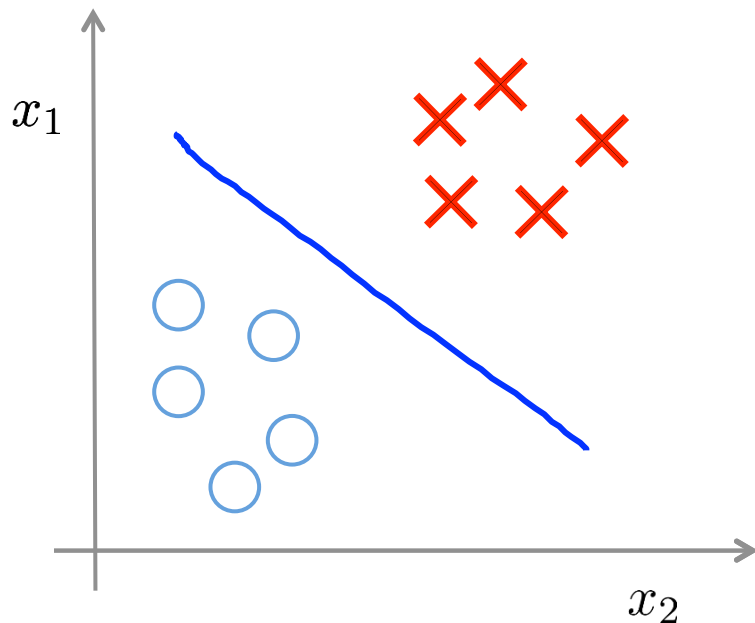# Clustering

## Unsupervised learning introduction

Machine Learning

# Supervised learning



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

# Unsupervised learning



$x_1$

Clustering algorithm

$x_2$

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

# Applications of clustering



Market segmentation



Social network analysis



Organize computing clusters



Image credit: NASA/JPL-Caltech/E. Churchwell (Univ. of Wisconsin, Madison)
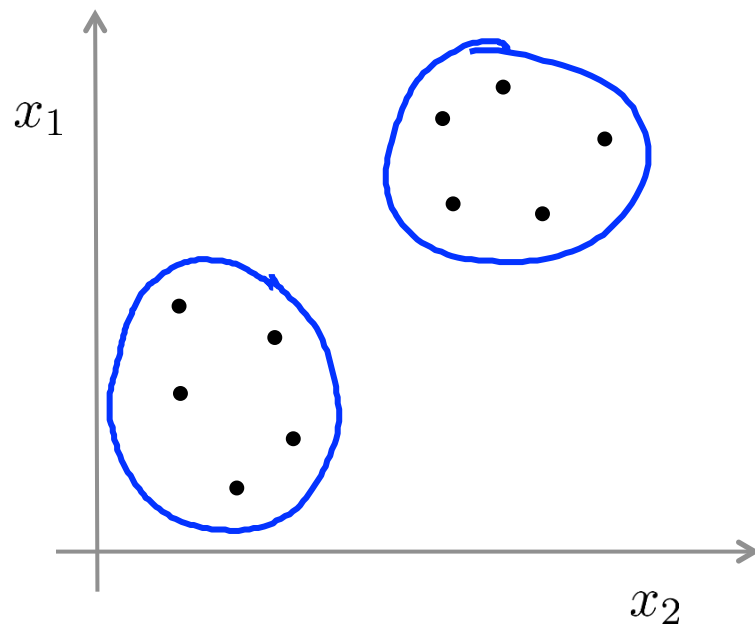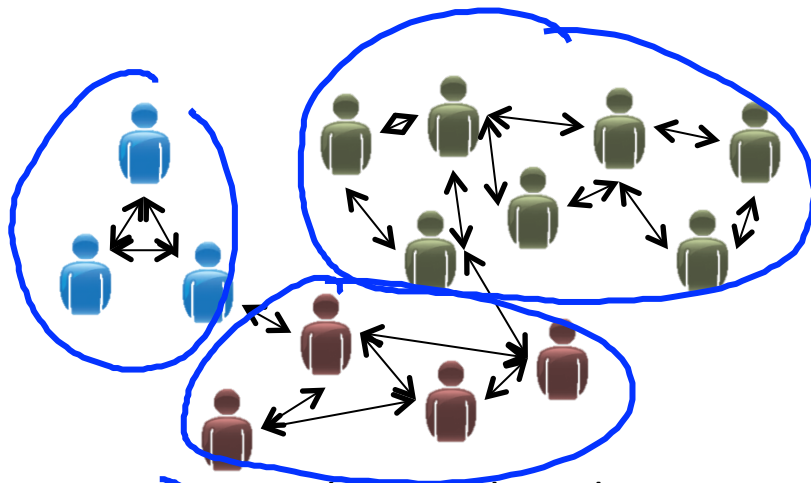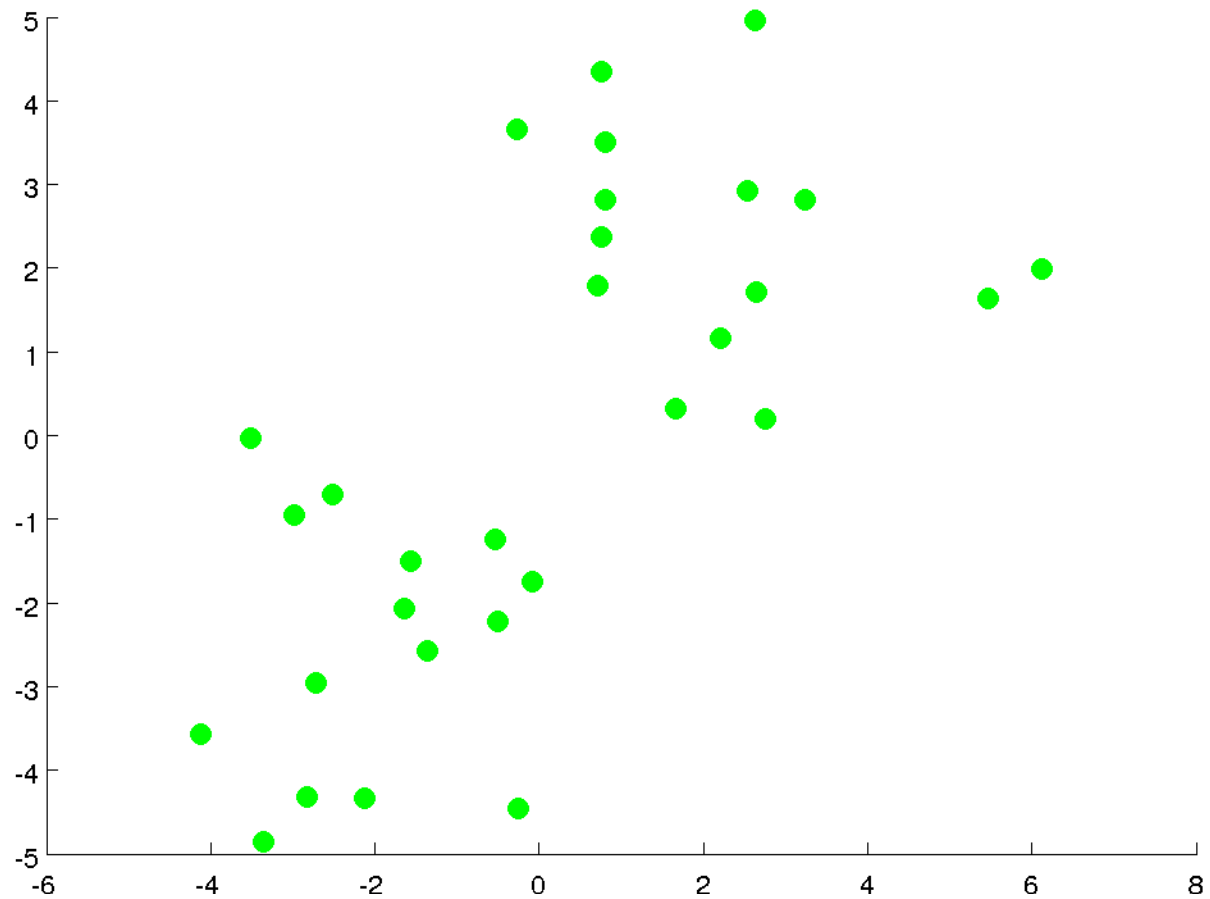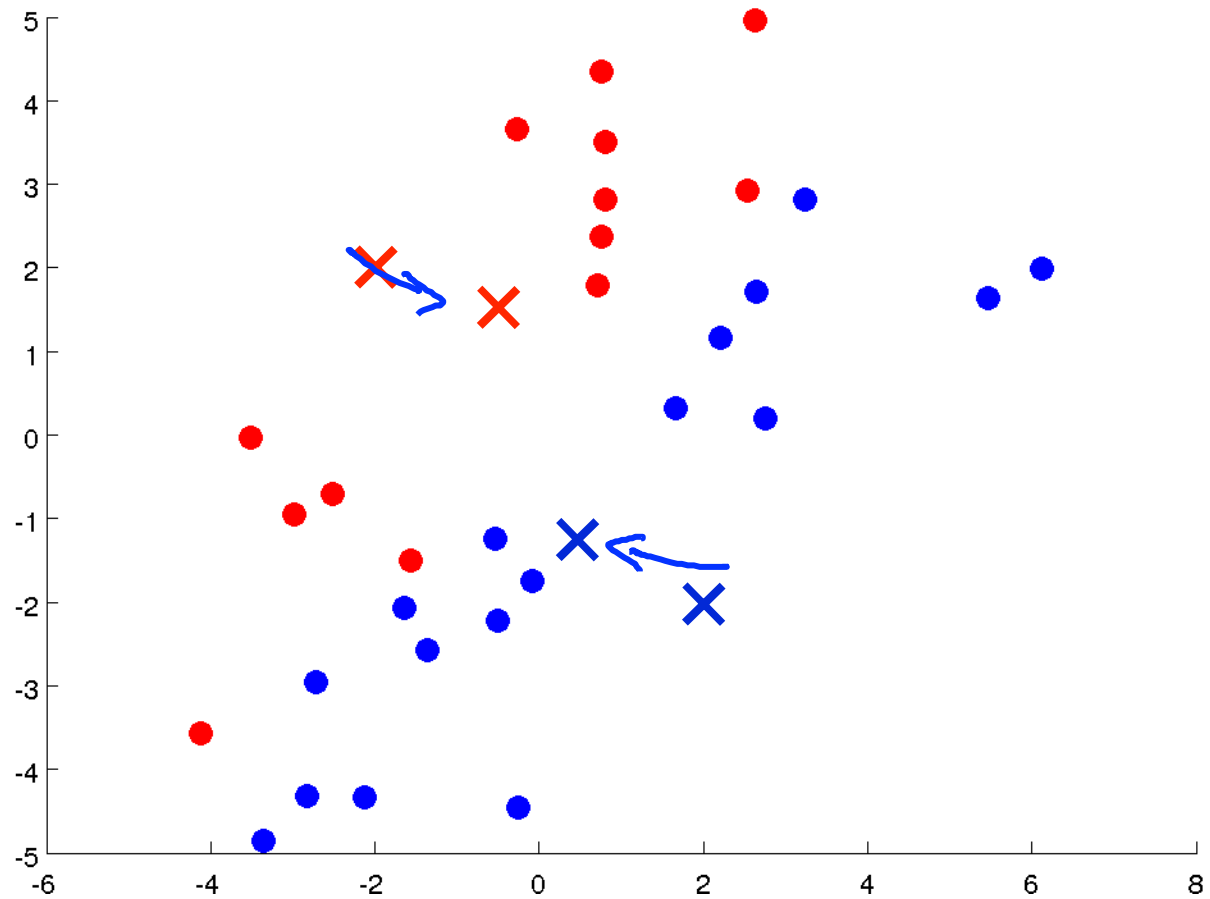
Astronomical data analysis

Andrew Ng

# Clustering

# K-means algorithm

Machine Learning

cluster centroids

Andrew Ng

# 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

$\mu_1$ ×   $\mu_2$ ×

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

Repeat {

*Cluster assignment step*

   for $i$ = 1 to $m$

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

   $\min_k \| x^{(i)} - \mu_k \|^2$

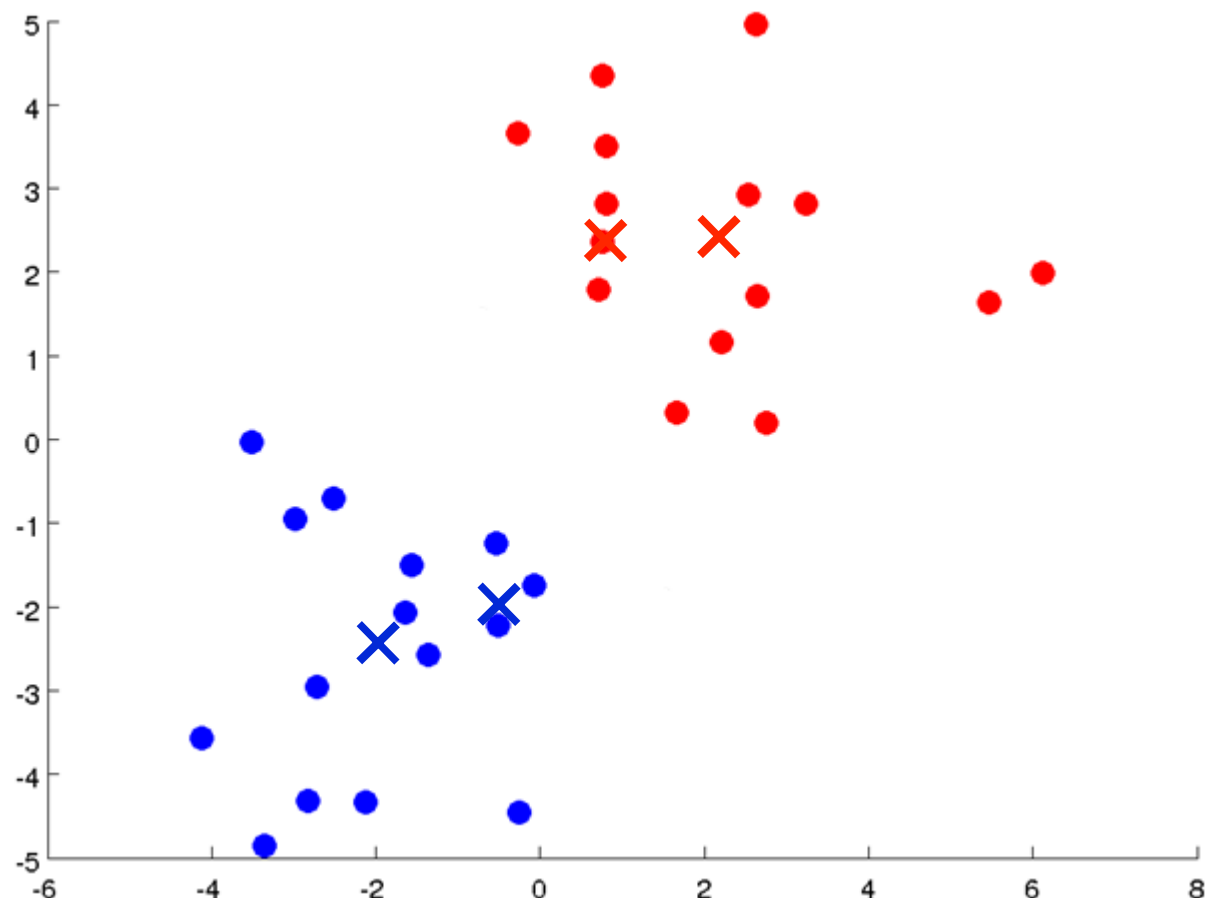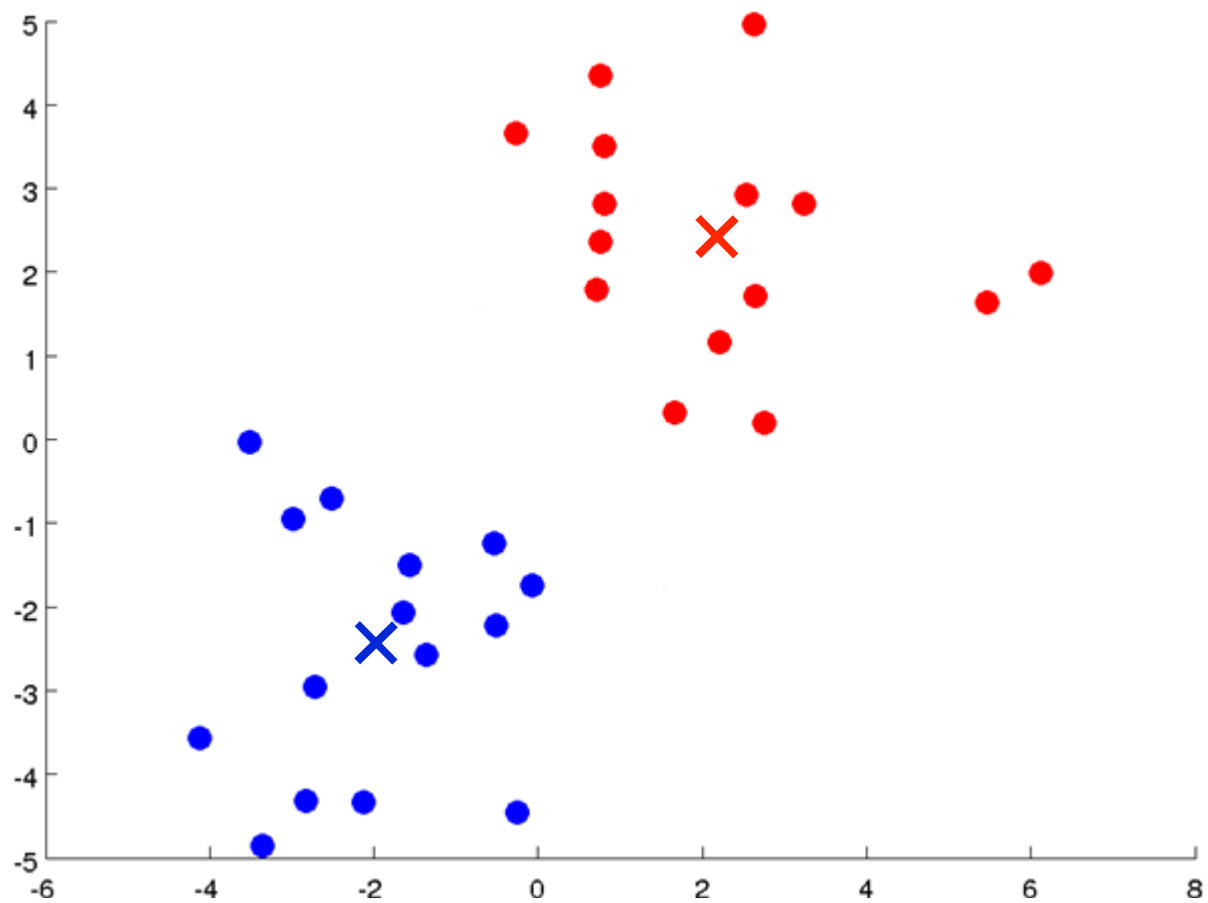   $\hookrightarrow c^{(i)}$

*More centroid*

   for $k$ = 1 to $K$

   $\mu_k$ := average (mean) of points assigned to cluster $k$
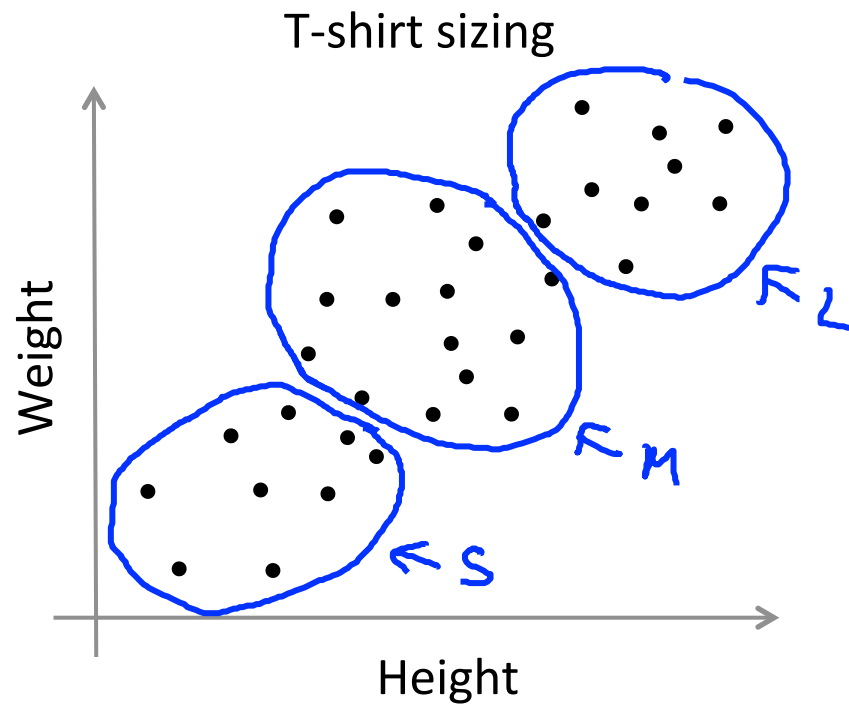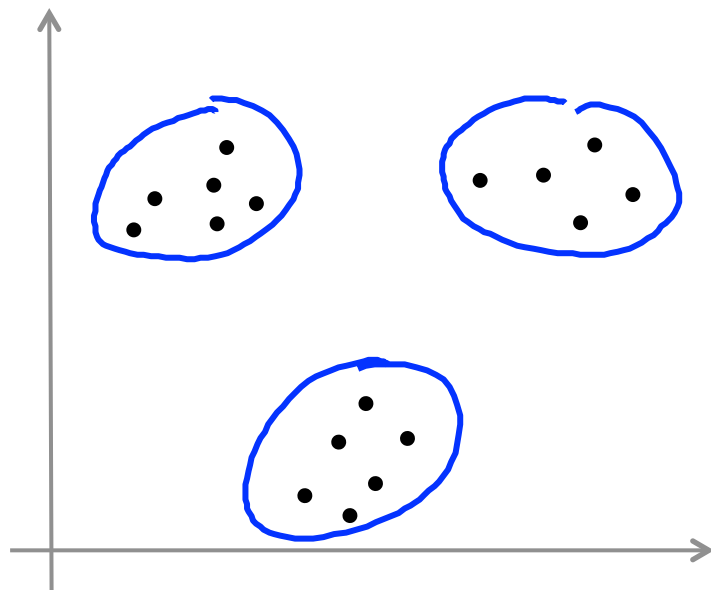
   $x^{(1)}, x^{(5)}, x^{(6)}, x^{(10)}$  $\rightarrow c^{(1)}=2, \quad c^{(5)}=2, \quad c^{(6)}=2,$
       $c^{(10)}=2$

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

}

# K-means for <u>non-separated clusters</u>

S, M, L

T-shirt sizing



Weight

Height

← S

← M

← L

Andrew Ng

Machine Learning

Clustering

Optimization objective

## K-means optimization objective

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

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

$K$ $\qquad$ $k \in \{1,2,...,k\}$

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

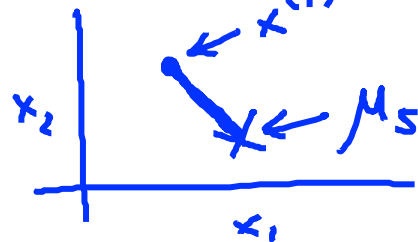$x^{(i)} \rightarrow 5 \qquad c^{(i)} = 5 \qquad \mu_{c^{(i)}} = \mu_5$

Optimization objective:

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

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

Distortion

$x_2$ $\qquad x^{(i)} \qquad \mu_5$

$x_1$

# K-means algorithm

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

Cluster assignment step

Minimize $J(\ldots)$ wrt $c^{(1)}, c^{(2)}, \ldots, c^{(m)}$ ←

(holding $\mu_1, \ldots, \mu_k$ fixed)

Repeat {

   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 $k$

}

Minimize $J(\ldots)$ wrt $\mu_1, \ldots, \mu_k$

Andrew Ng

# Clustering

## Random initialization

Machine Learning

**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 $k$
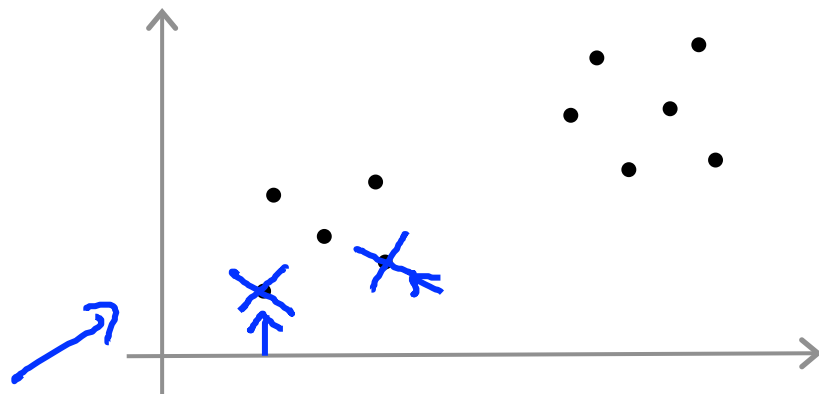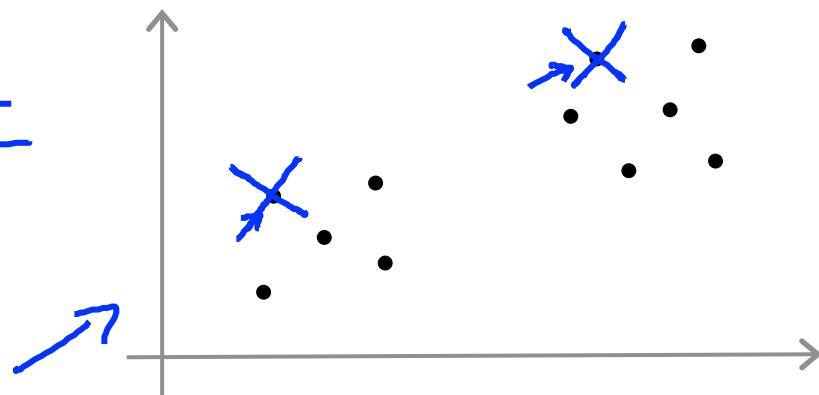
}

# Random initialization
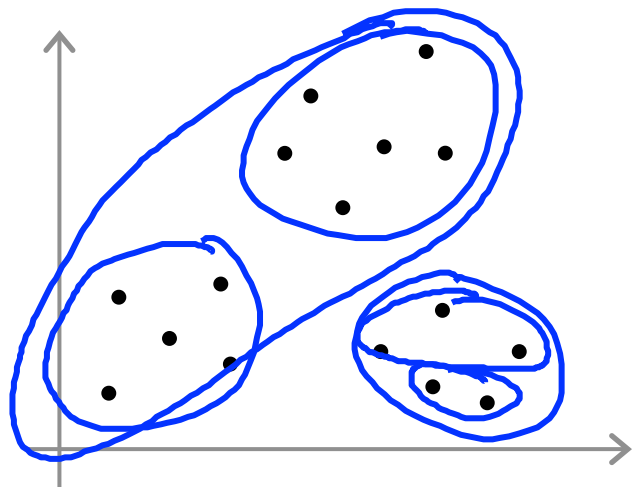
Should have $K < m$

Randomly pick $K$ training examples.

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



K=2

$\mu_1 = x^{(i)}$

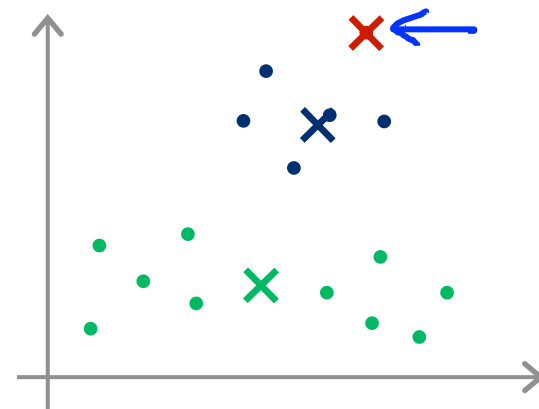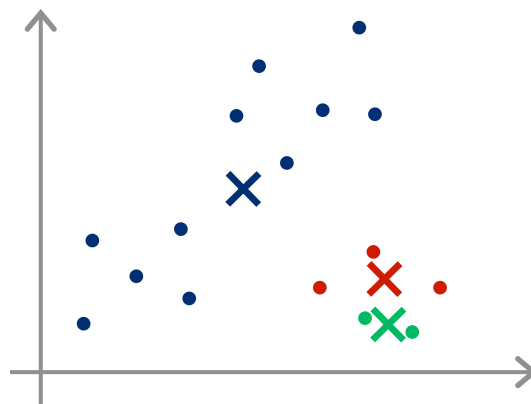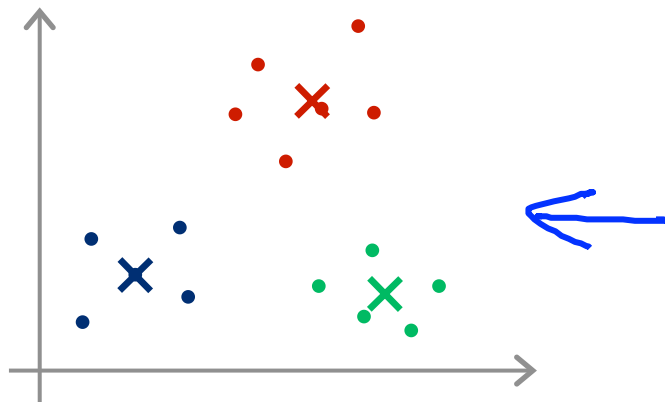$\mu_2 = x^{(j)}$

# Local optima



$$J\left(c^{(1)}, \ldots, c^{(m)}, \mu_1, \ldots, \mu_k\right)$$

**Random initialization**

For i = 1 to 100 {

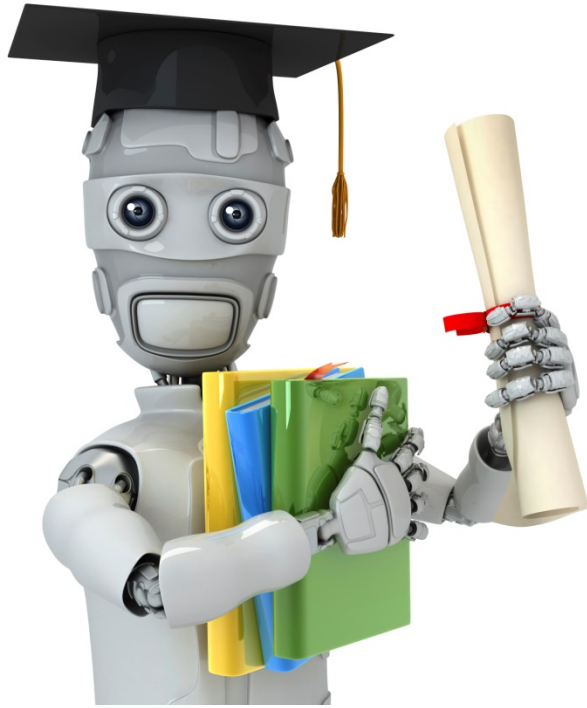        Randomly initialize K-means.

        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 $J(c^{(1)}, \ldots, c^{(m)}, \mu_1, \ldots, \mu_K)$
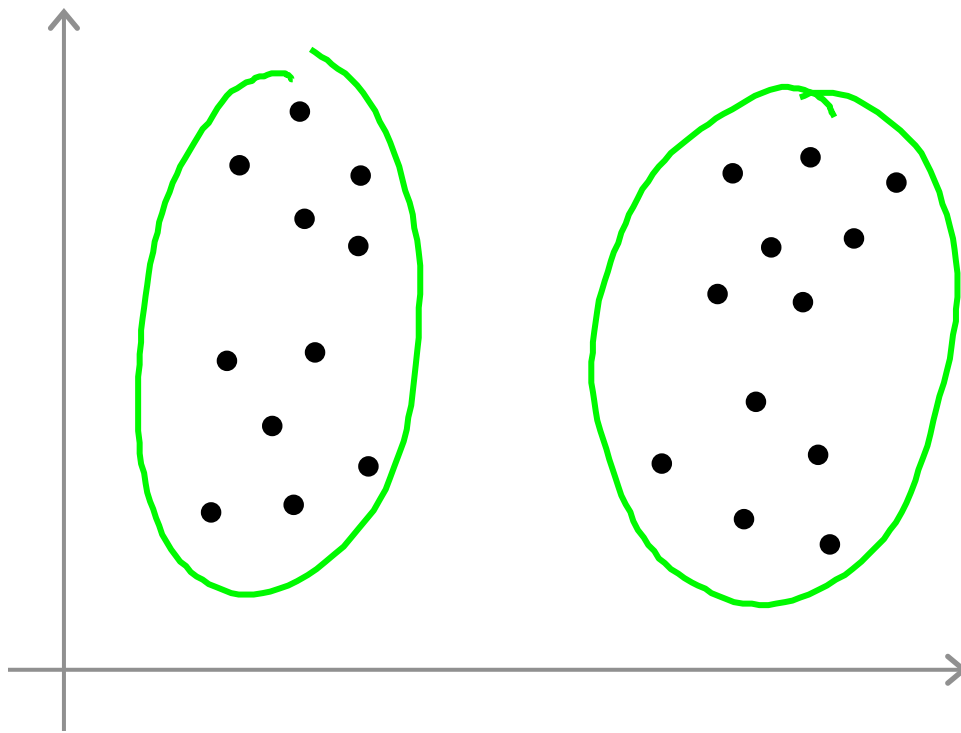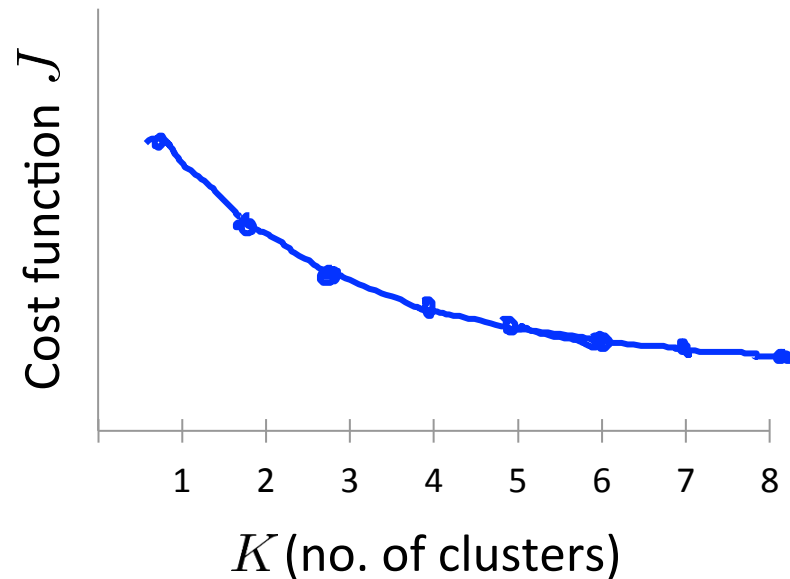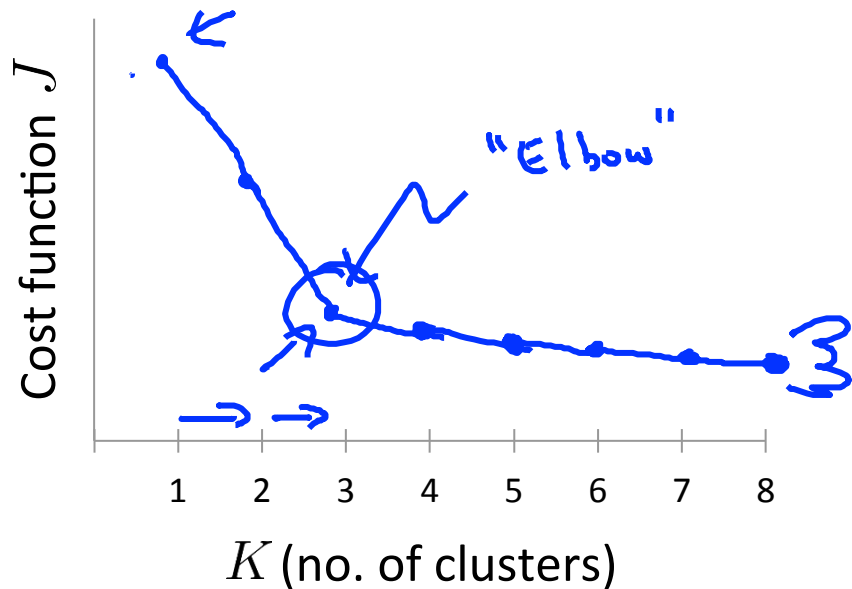
# Clustering

## Choosing the number of clusters

Machine Learning

# What is the right value of K?

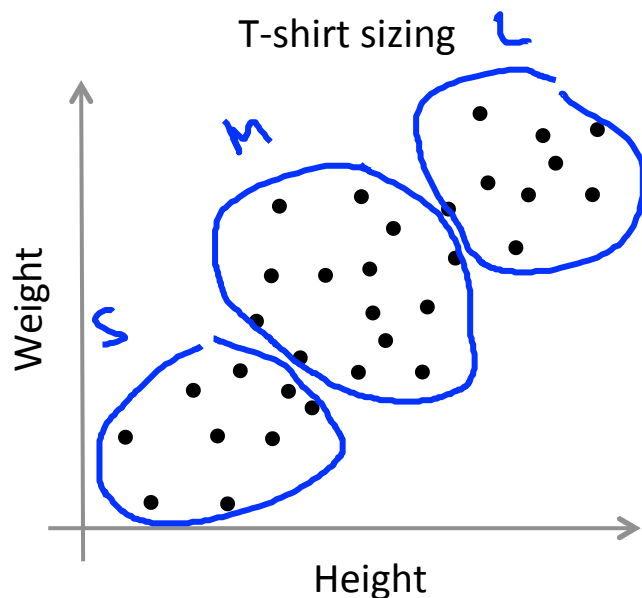# Choosing the value of K

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.

K=3    S, M, L

K=5    XS, S, M, L, XL

E.g.



T-shirt sizing

L

M

S

Weight

Height



T-shirt sizing

XL

L

M

S

XS

Weight

Height