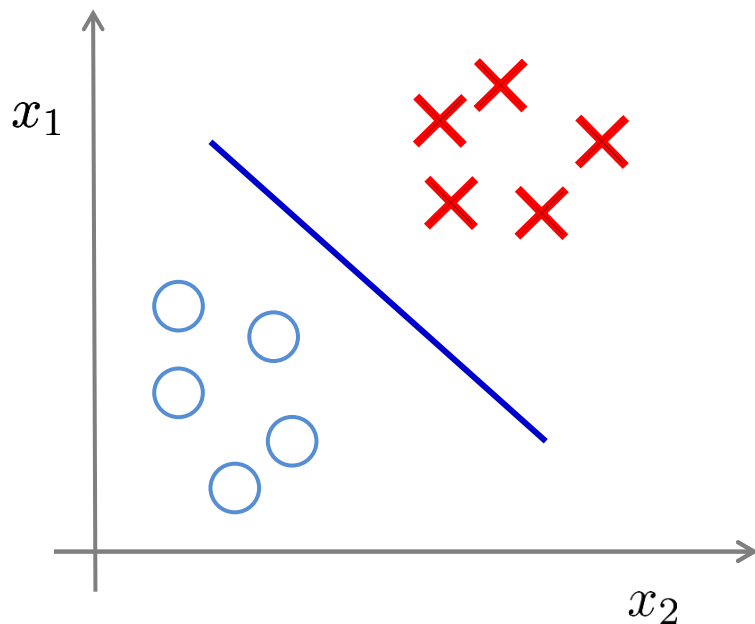# Clustering

## Unsupervised learning introduction

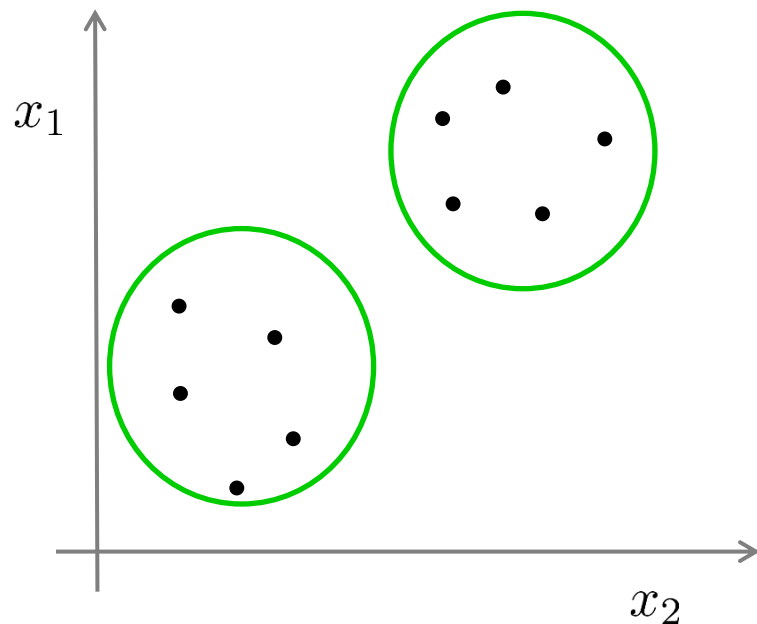Machine Learning

# Supervised learning



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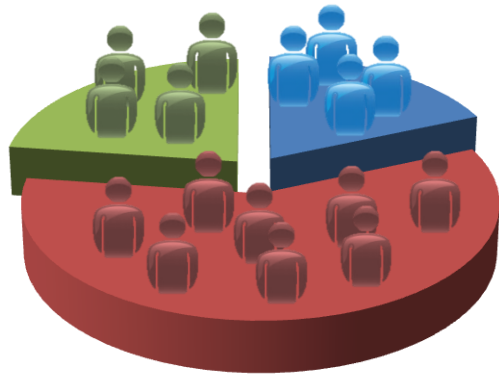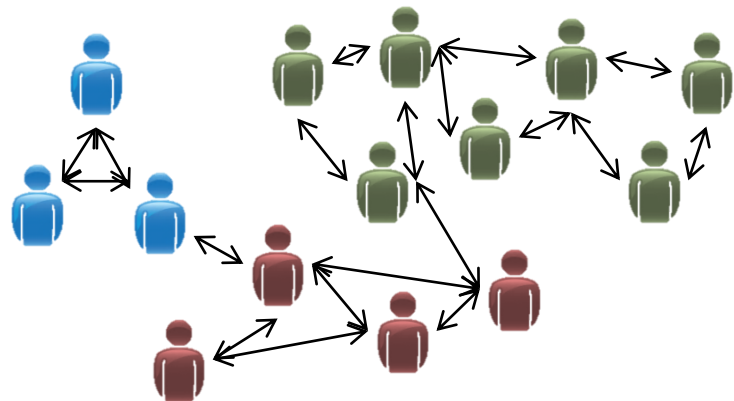
# Unsupervised learning



Clustering algorithm

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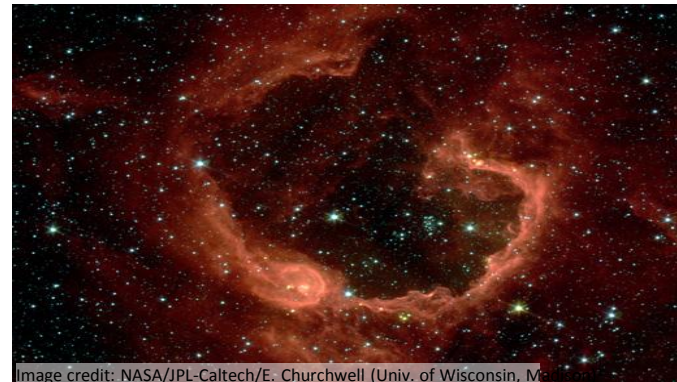
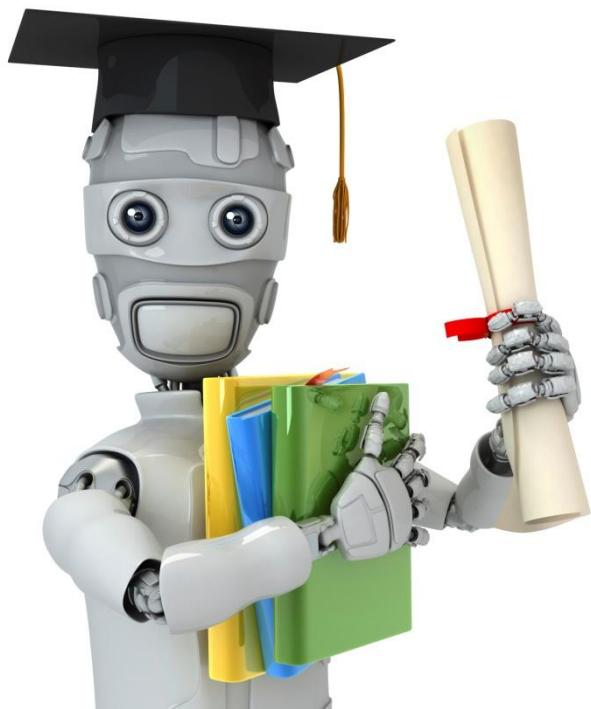# Applications of clustering



Market segmentation

Social network analysis
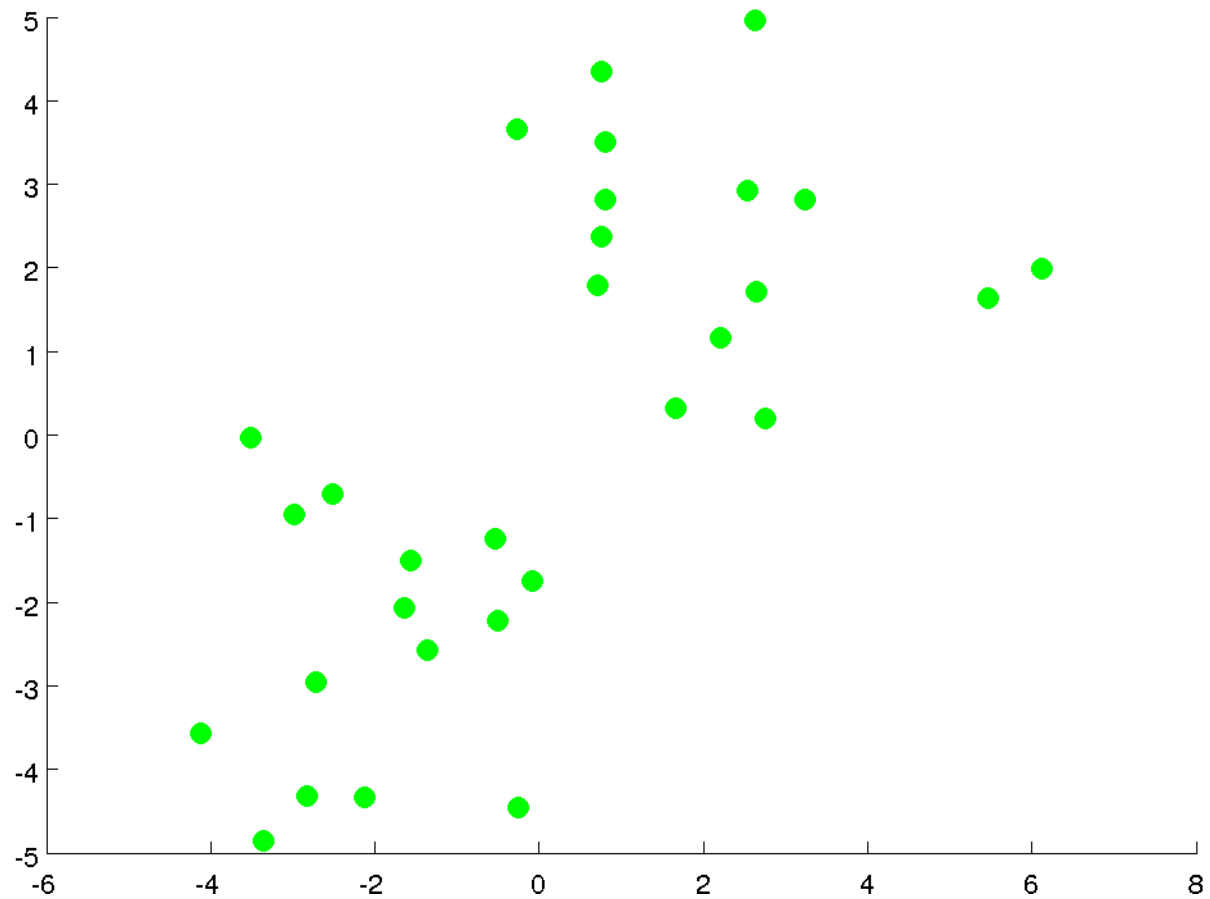
Organize computing clusters

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

Astronomical data analysis
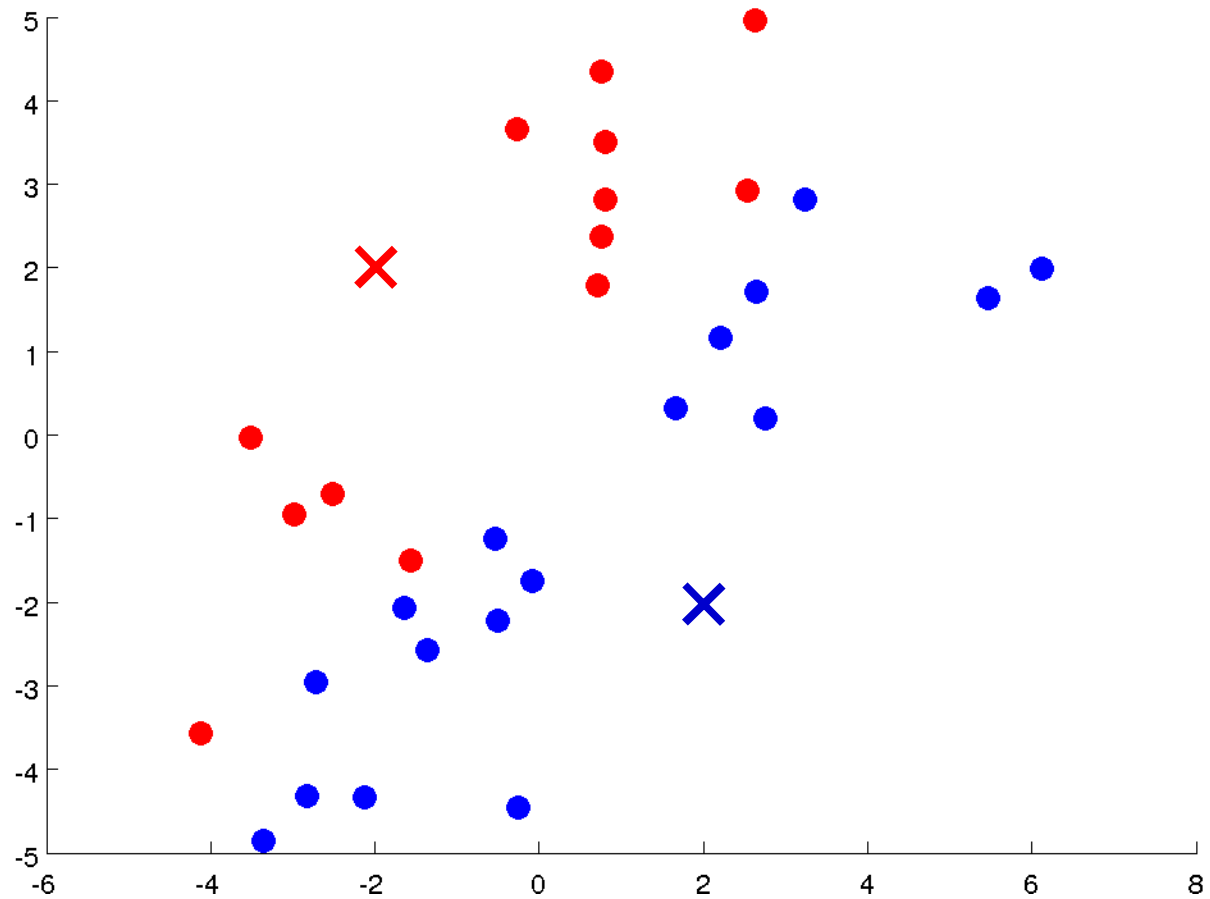
Andrew Ng

# Clustering

# K-means algorithm

Machine Learning

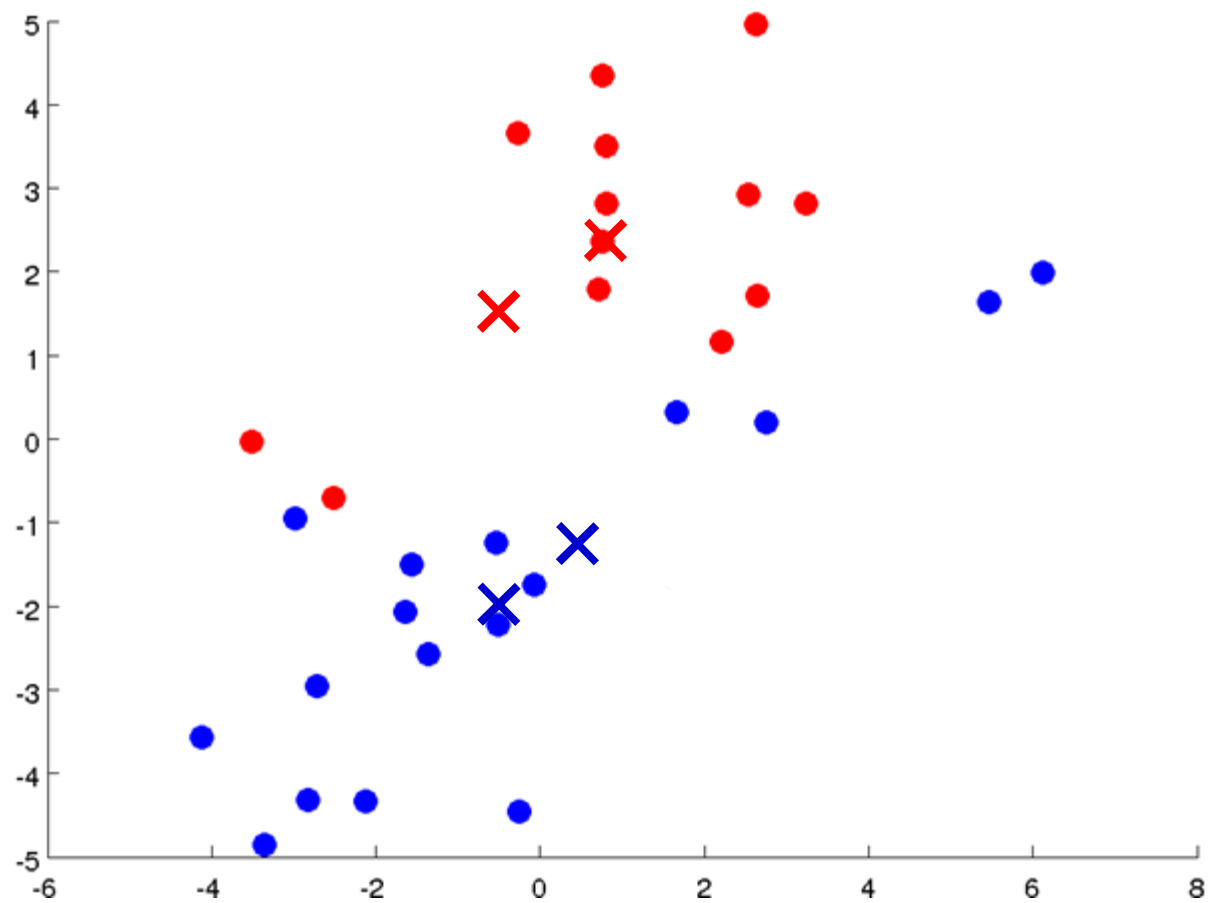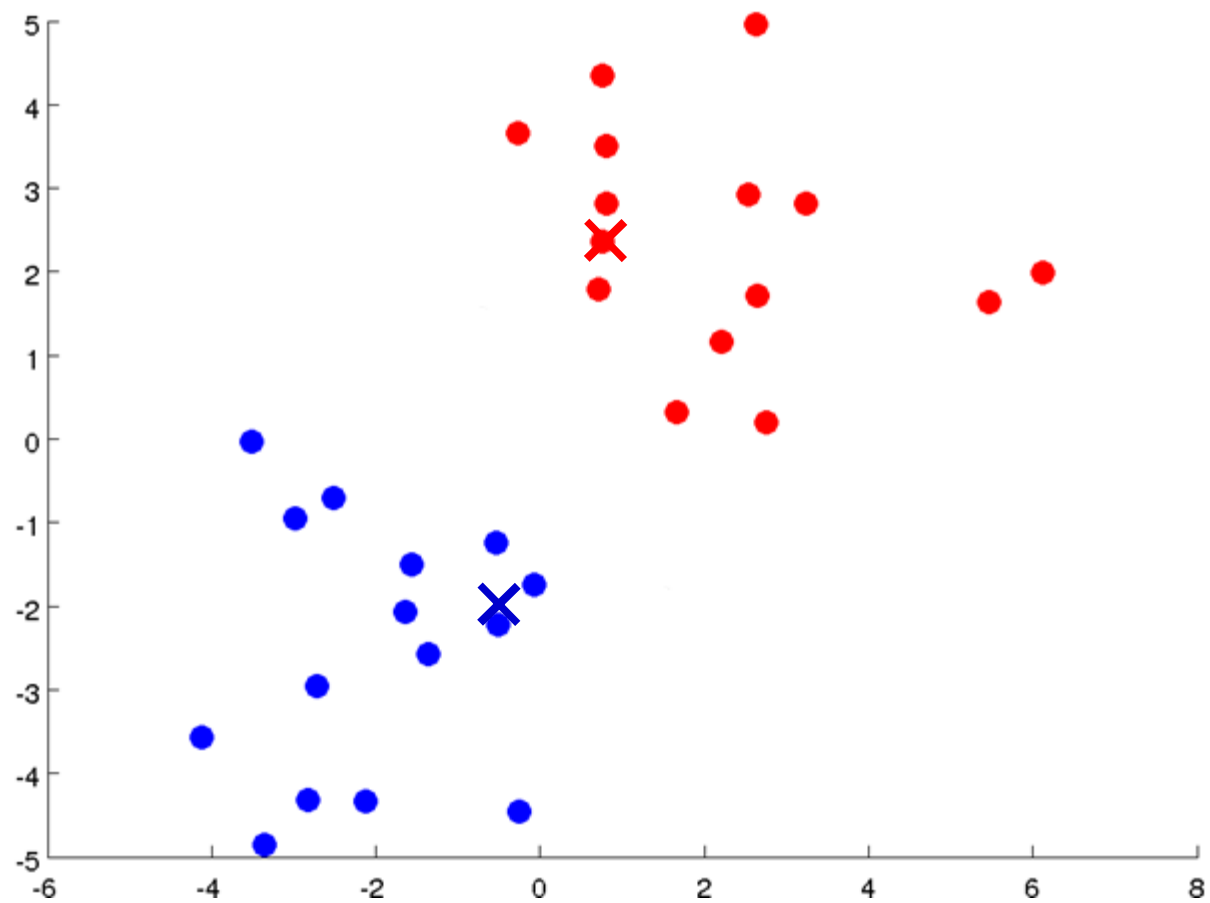Cluster centroids

Andrew Ng

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

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 $x^{(i)}$ $\quad \boldsymbol{\min_{k}} = ||\boldsymbol{x}^{(i)} - \boldsymbol{u_k}||^2$

Move centroids step

for $k$ = 1 to $K$

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

}

# K-means for non-separated clusters



separable

T-shirt sizing

S,M,L

Weight

Height

S

M

L

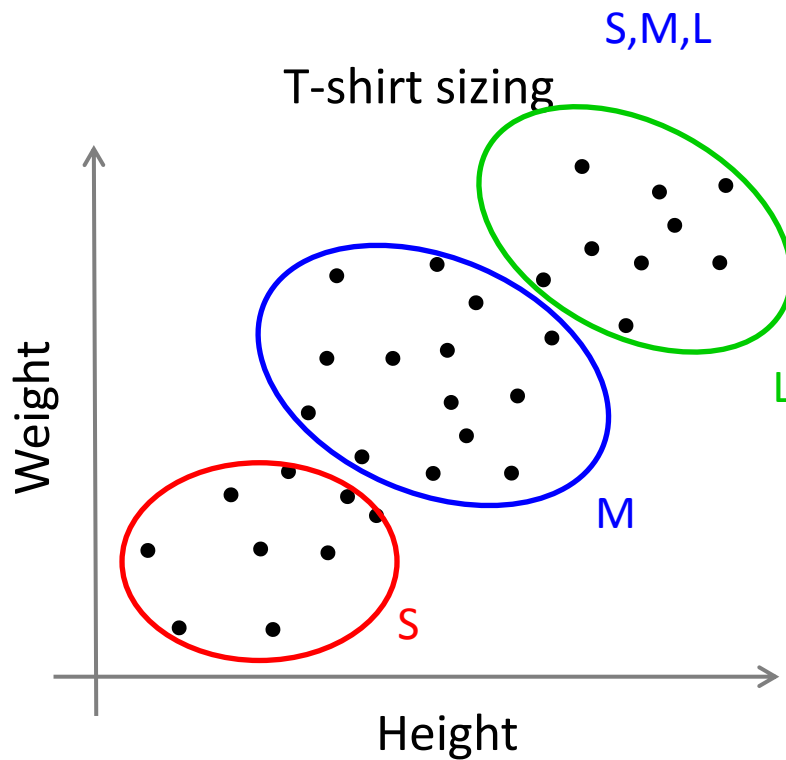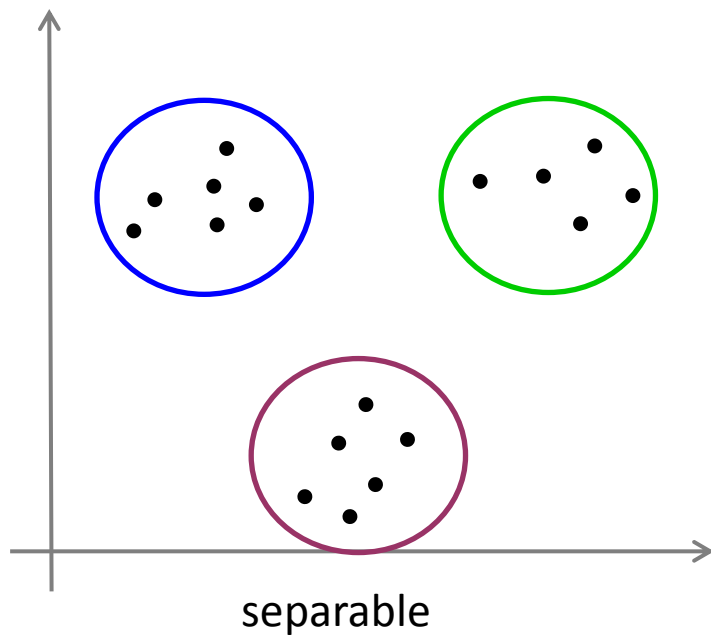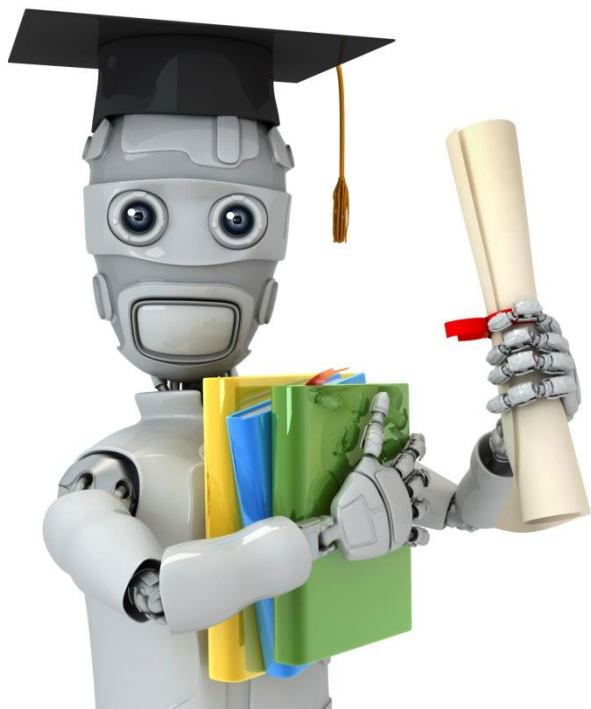Andrew Ng

Clustering

Optimization objective

Machine Learning

# 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

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

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 {

Cluster assignment step

$Minimize\ J(\ldots)\ w.r.t\ c^{(1)}, c^{(2)}, \ldots, c^{(k)}$
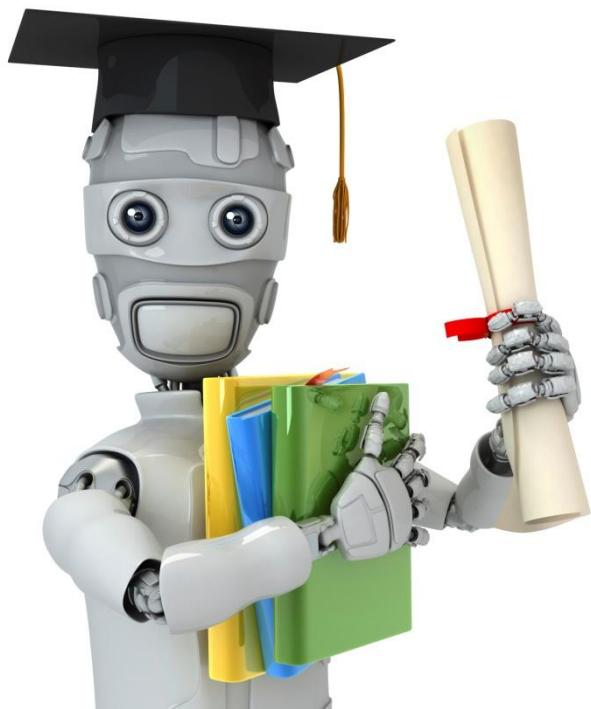$(holding\ \mu_1, \ldots, \mu_k\ fixed)$

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 $k$

}

Move centroids step        $Minimize\ J(\ldots)\ w.r.t\ \mu_1, \ldots, \mu_k$

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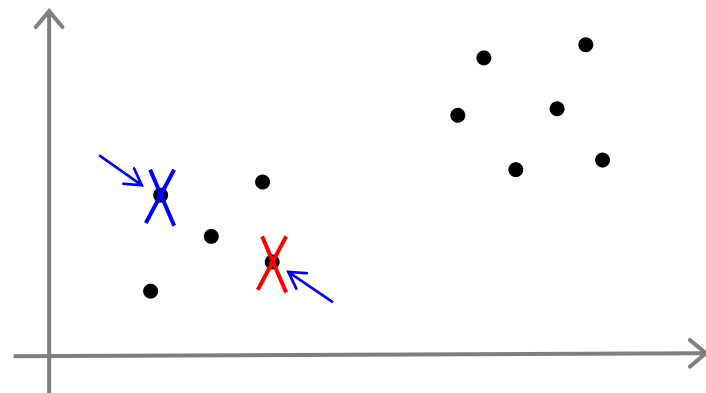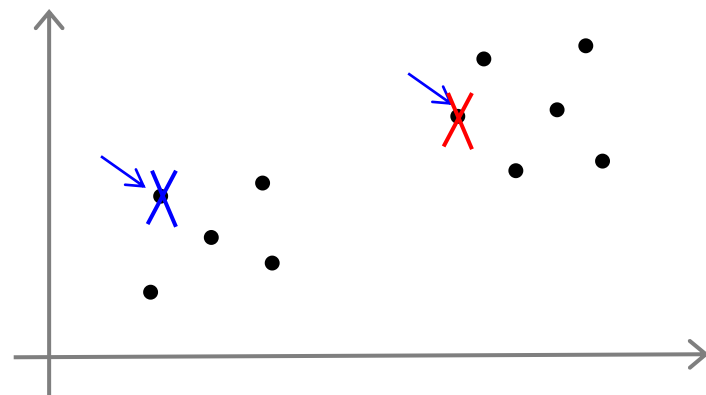           $\mu_k$ := average (mean) of points assigned to cluster $k$
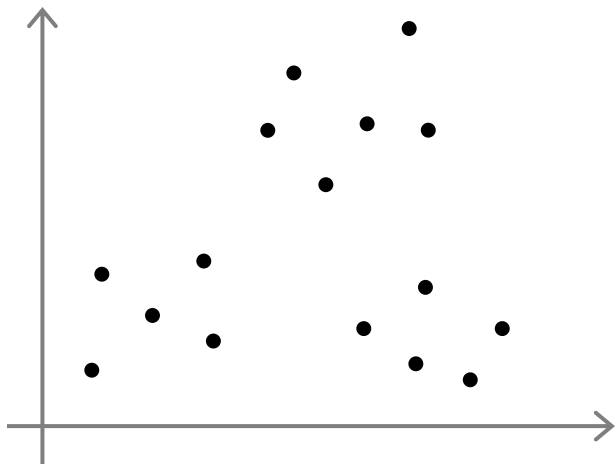       }

# Random initialization

Should have $K < m$

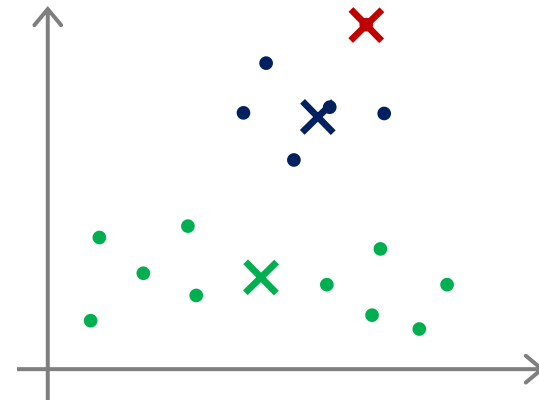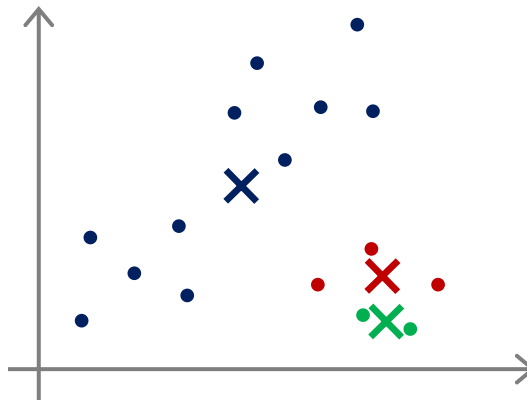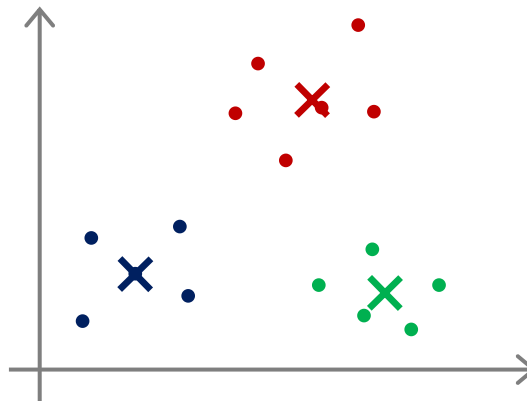Randomly pick $K$ training examples.

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

**Local optima**

But.....
Unlucky random
initialization

**Random initialization**

For i = 1 to 100 {

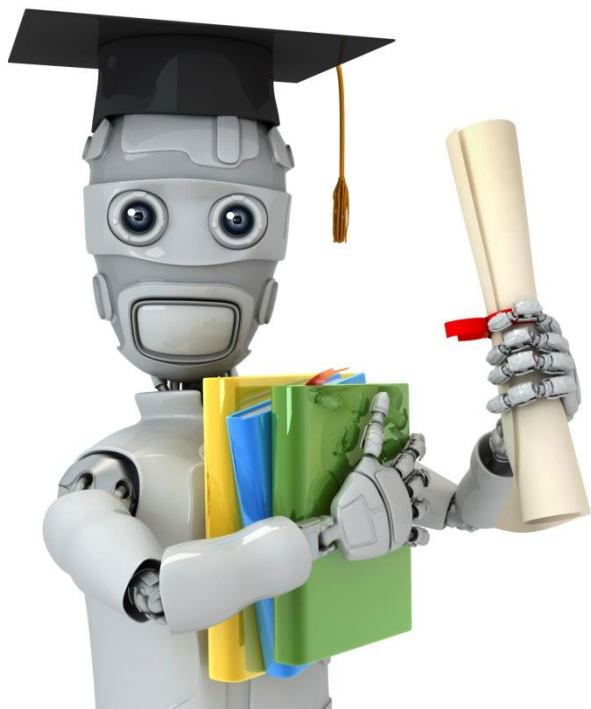       Randomly initialize K-means.
       Run K-means. Get $c^{(1)}, \ldots, c^{(m)}, \mu_1, \ldots, \mu_K.$
       Compute cost function (distortion)
$$J(c^{(1)}, \ldots, c^{(m)}, \mu_1, \ldots, \mu_K)$$
       }

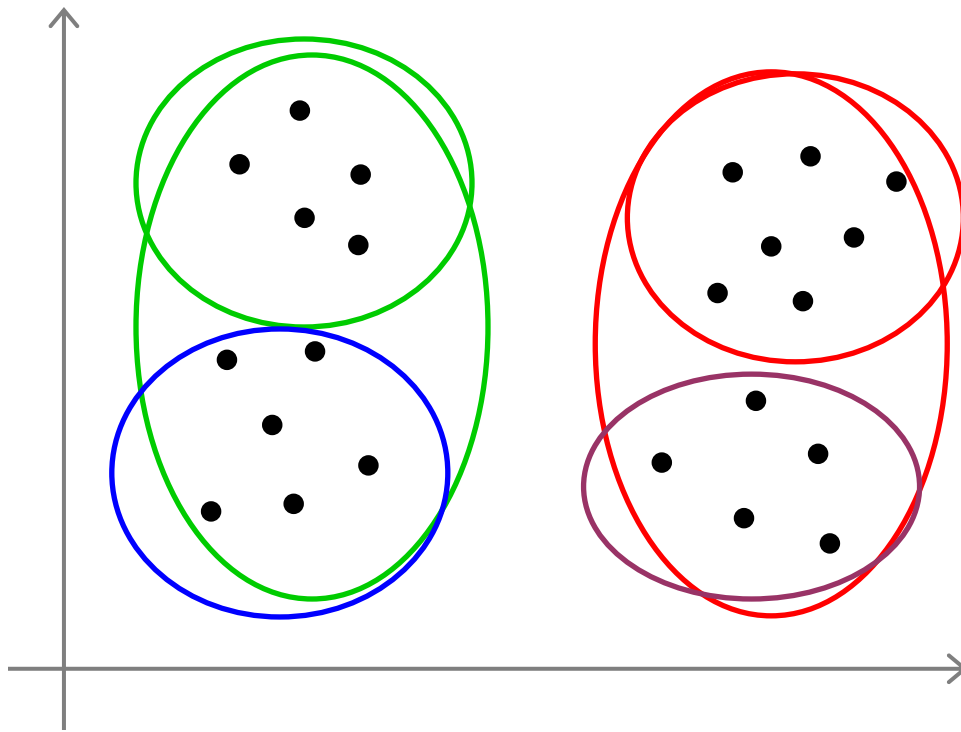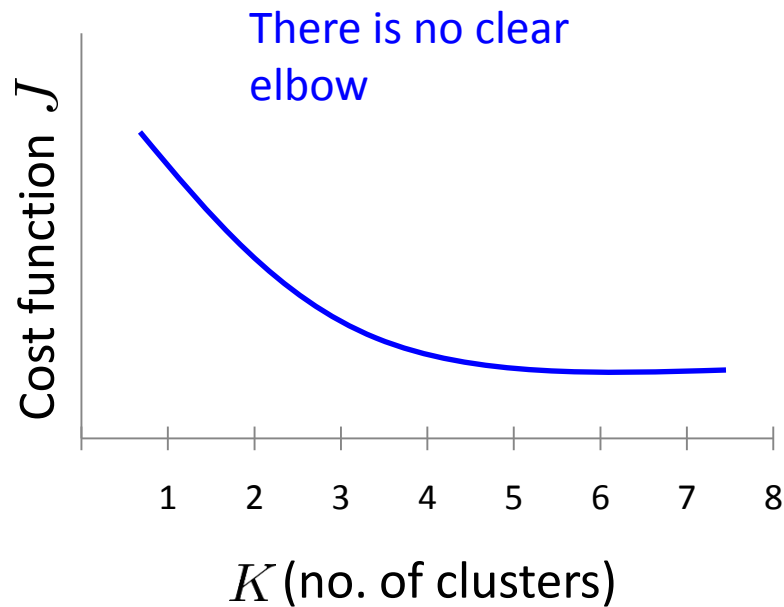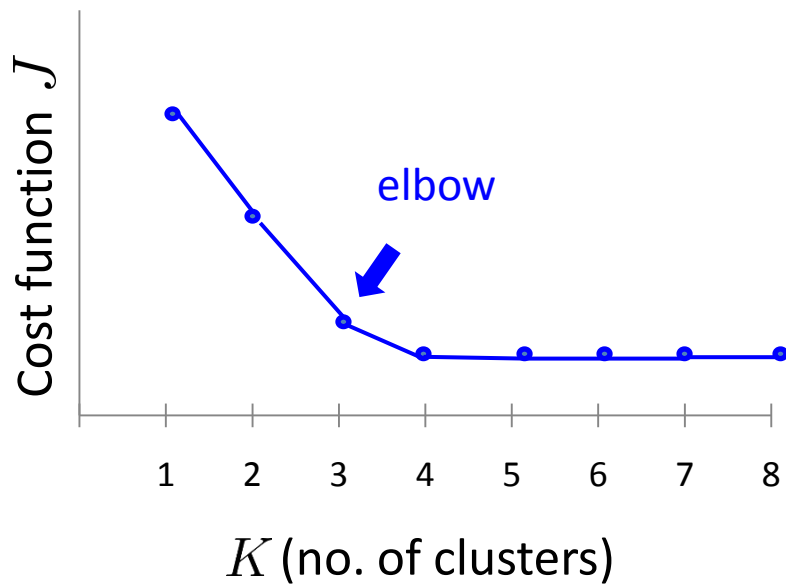Pick clustering that gave lowest cost $J(c^{(1)}, \ldots, c^{(m)}, \mu_1, \ldots, \mu_K)$

# Clustering

## Choosing the number of clusters

Machine Learning

# What is the right value of K?

# Choosing the value of K

Elbow method:



Cost function $J$

elbow

$K$ (no. of clusters)

Cost function $J$

There is no clear elbow

$K$ (no. of clusters)

# 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.