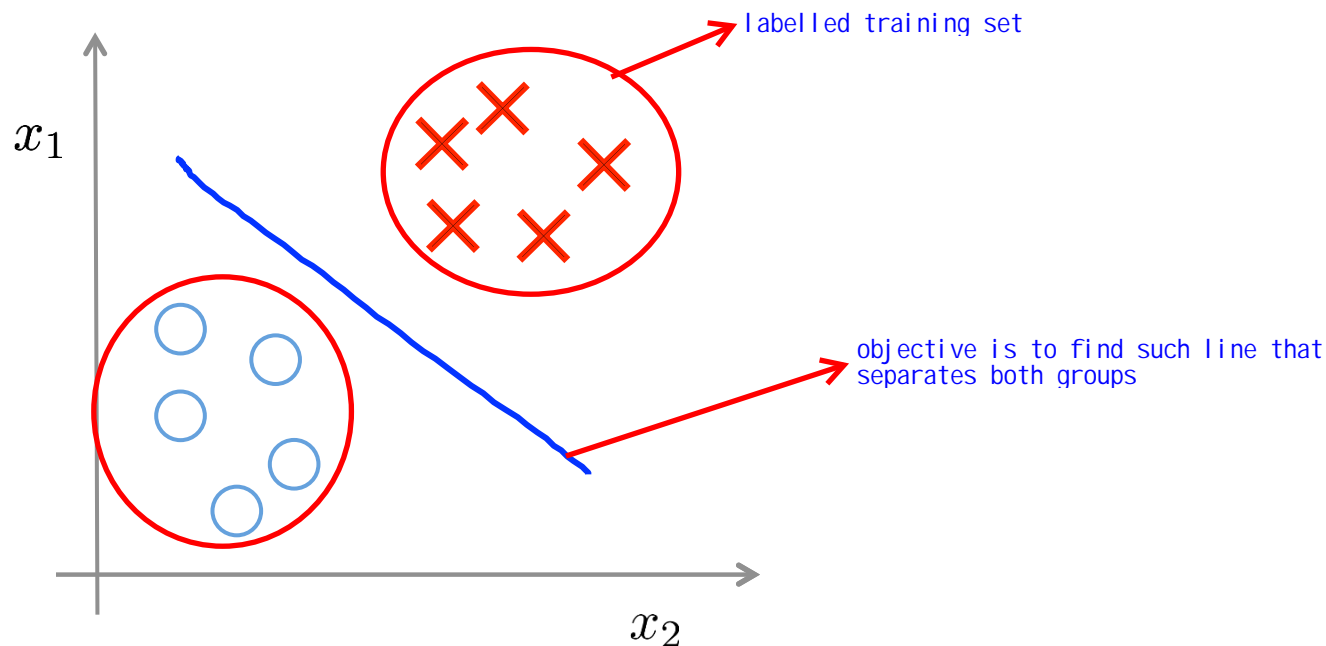# Clustering

## Unsupervised learning introduction

Learn From unlabelled data!

Machine Learning
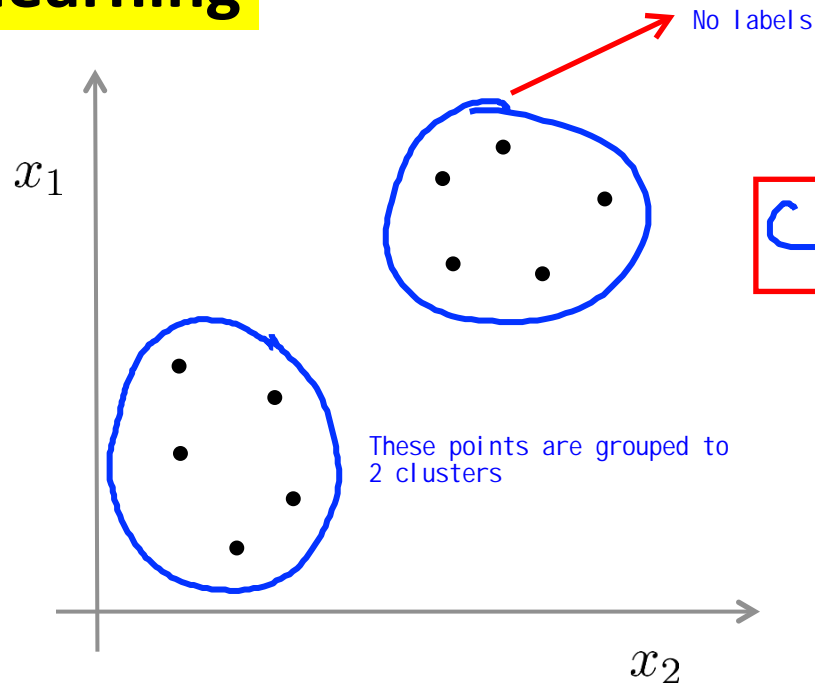
# Supervised learning



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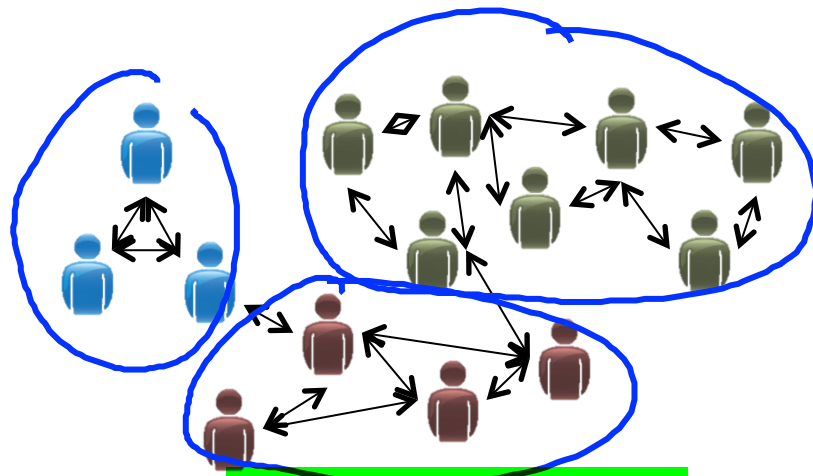
# Unsupervised learning



No labels

$x_1$

Clustering algorithm

These points are grouped to 2 clusters

$x_2$

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

# Applications of clustering


Market segmentation


Social network analysis



In data center, some computers tend to work together.
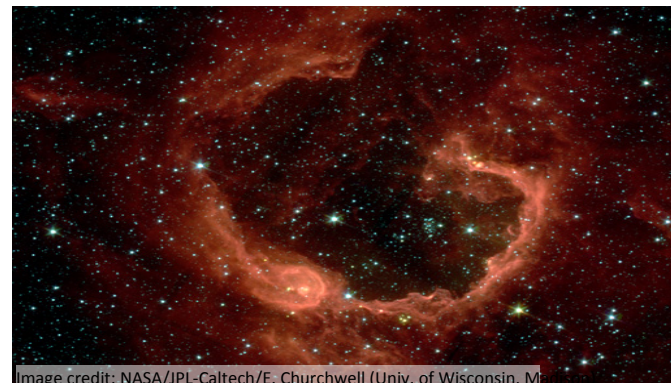
Organize computing clusters


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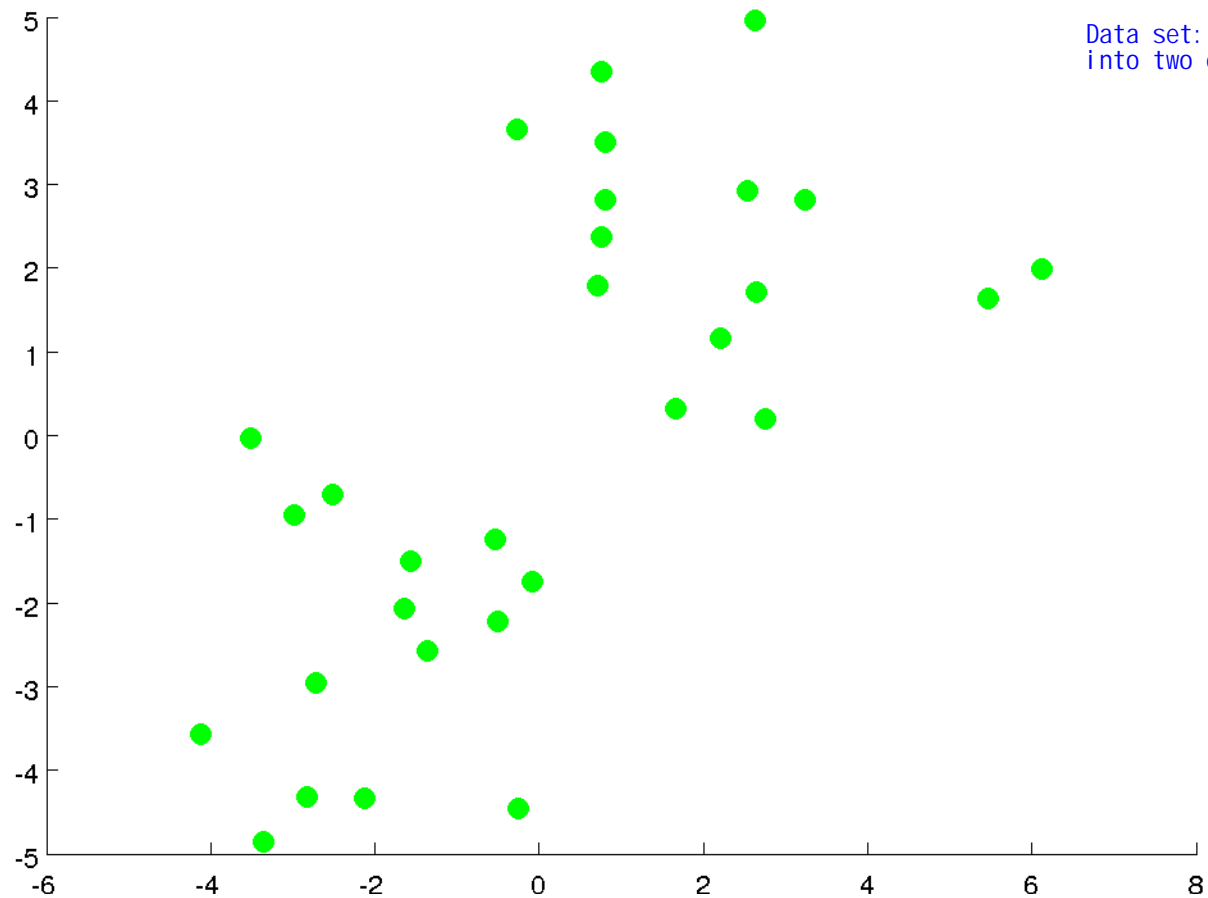
Astronomical data analysis

Andrew Ng

Machine Learning

# Clustering

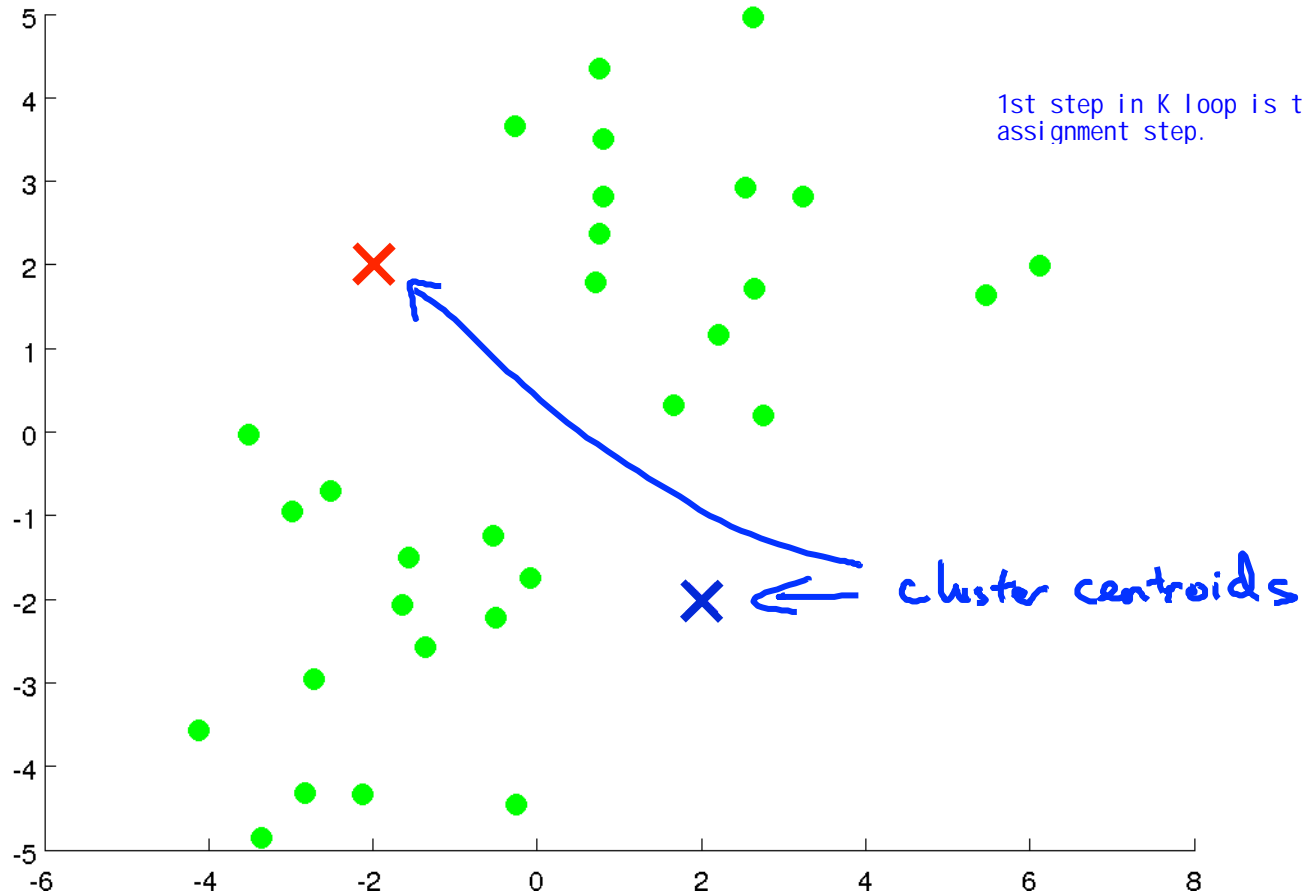## K-means algorithm

iterative algorithm

By far, the most popular and most commonly used algorithm

1st step: randomly initialize two points called cluster centroids.

1st step in K loop is the cluster assignment step.

cluster centroids

Andrew Ng

go through all the data points and label them by either red or blue depending on whether the data point is closer to the red centroid or the blue centroid.

The 2nd step in the K means loop is the moving centroid step.

we look at all the red points and compute the mean of all the red points and we then move our centroid there. We perform the same for blue points.

Andrew Ng

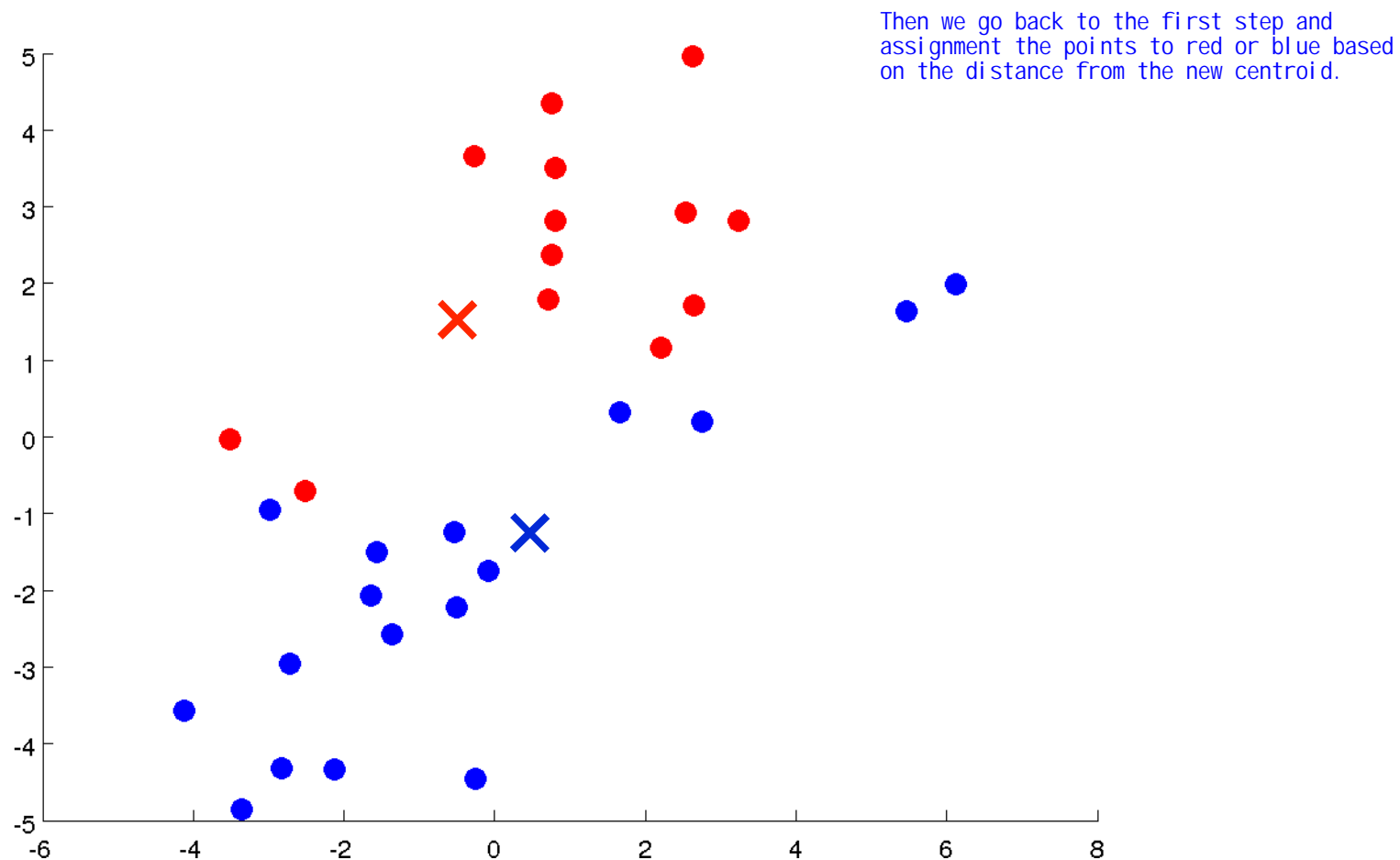Then we go back to the first step and assignment the points to red or blue based on the distance from the new centroid.

Andrew Ng

Andrew Ng

keep running

Andrew Ng

keep running, eventually, the locations of the centroids will not change further.

At this stage, we say that the K means algorithm has converged.

Andrew Ng

# K-means algorithm

Input:

For now, we just decide certain number for K.

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

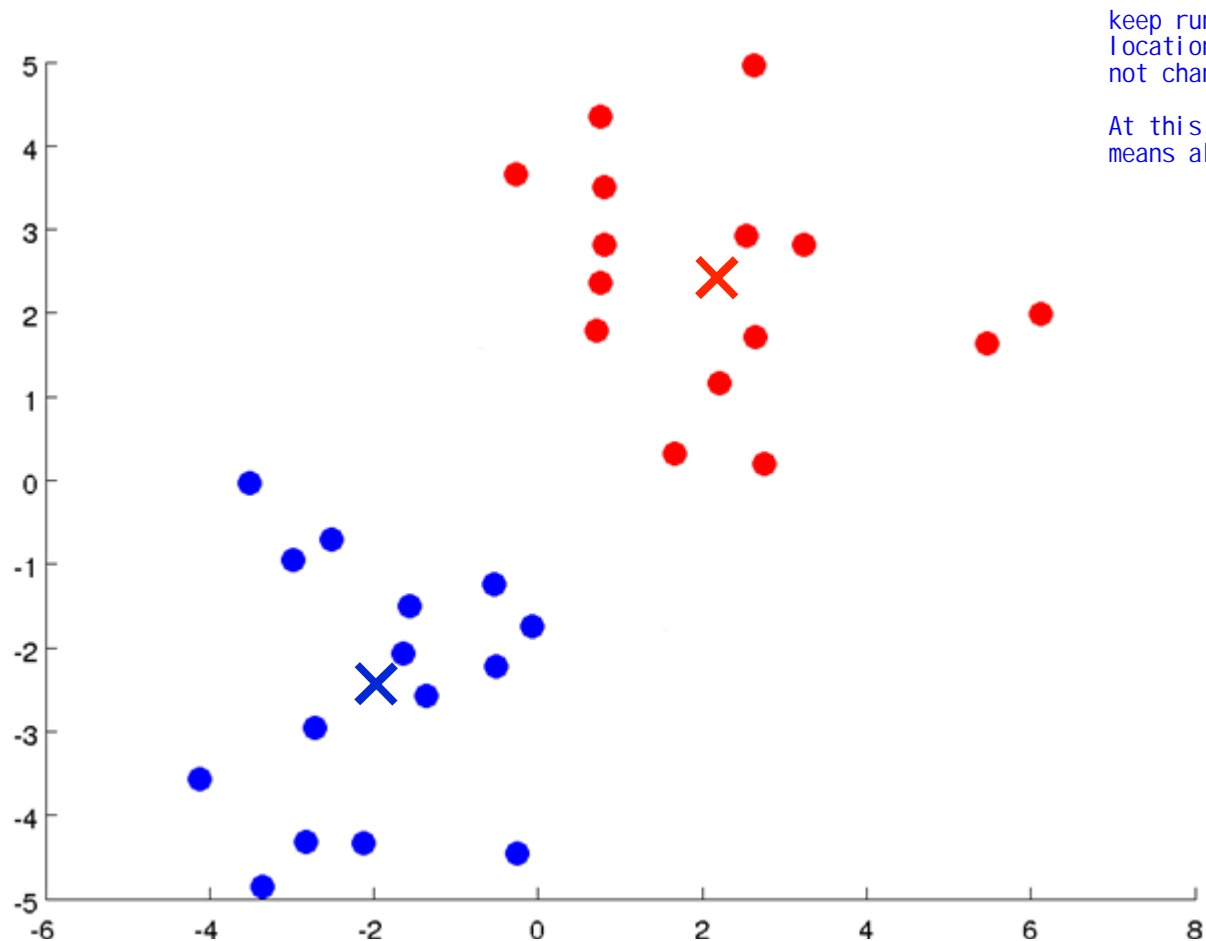$x^{(i)} \in \mathbb{R}^n$ (drop $x_0 = 1$ convention)

Andrew Ng

# K-means algorithm

cluster assignment step

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

1st step

$\mu_1$ ✕   $\mu_2$ ✕

Repeat {

Cluster assignment step

   for $i$ = 1 to $m$

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

is a number from 1 to K, which indicates which centroid it is closet to x_i

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

$\to c^{(i)}$

Move centroid

   for $k$ = 1 to $K$

move centroid step

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

$x^{(1)}, x^{(5)}, x^{(6)}, x^{(10)}$

$\to c^{(1)} = 2, \quad c^{(5)} = 2, \quad c^{(6)} = 2,$

$c^{(10)} = 2$

}

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

average

when there is cluster with no points: we can just eliminate that cluster then we have K-1 clusters

Andrew Ng

# K-means for non-separated clusters

One other common application of K means

S, M, L

An example of market segregation

T-shirt sizing



Weight

Height

well-separated

large

medium

small

not well-separated data set
we can also apply K-means

Andrew Ng

# Clustering

## Optimization objective

K means algorithm also has an optimization objective!

Machine Learning

## K-means optimization objective

we keep tracking these two parameters

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

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

$K$        $k \in \{1, 2, …, k\}$

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

$x^{(i)} \rightarrow 5$        $c^{(i)} = 5$        $\mu_{c^{(i)}} = \mu_5$
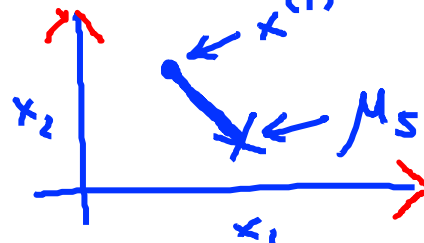
x_i is assigned to cluster 5

corresponding centroid location

## Optimization objective:

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

$x^{(i)}$

objective!

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

Distortion

$x_2$        $x^{(i)}$        $\mu_5$        $x_1$

Andrew Ng

# K-means algorithm

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

Repeat {

Cluster assignment step

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

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

obvious!

    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 $\mu_1, \ldots, \mu_k$

holding c_i

with respect to

Andrew Ng

# Clustering

## Random initialization

To make K-means avoid local optima

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 $k$

}

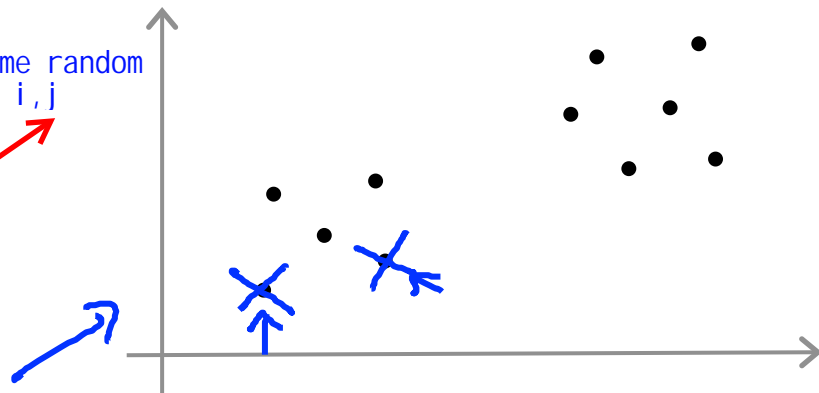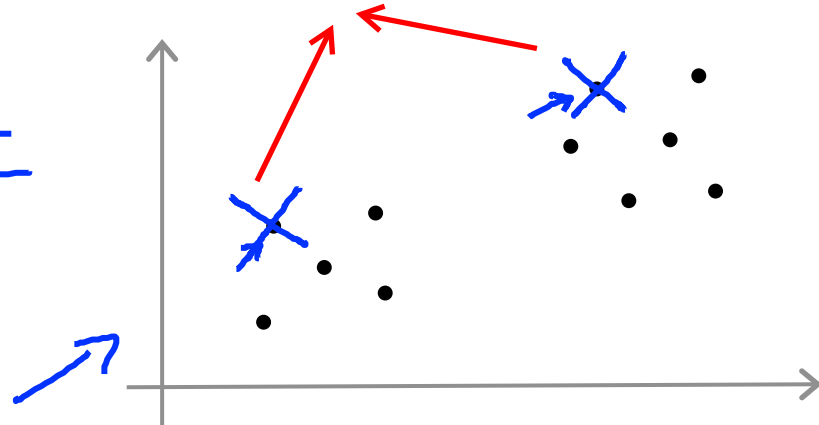# Random initialization

Should have $K < m$

Randomly pick $K$ training examples.

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

K-means can end up with different solution depending on your initial conditions

$$\mu_1 = x^{(i)}$$
$$\mu_2 = x^{(j)}$$
$$\vdots$$

$K = 2$

randomly pick two initial centroid

for some random values i,j
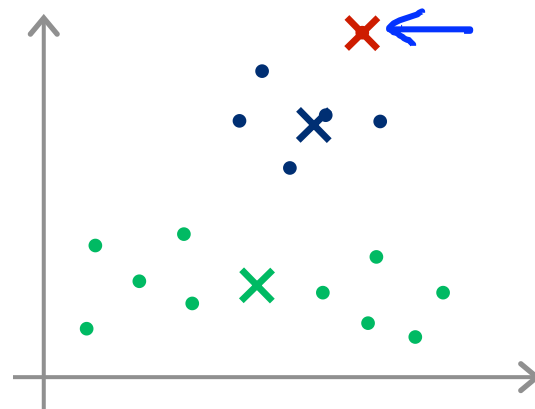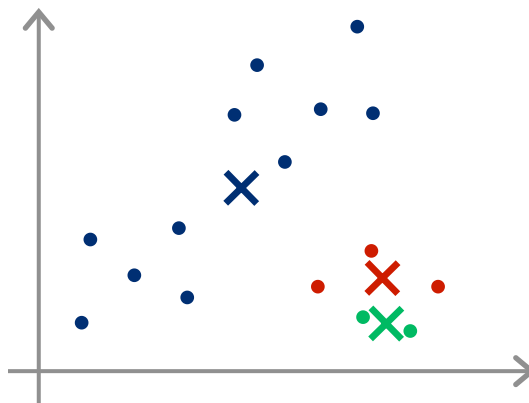
**Local optima**

Can end up at local optima
To avoid this, we can run K-means
multiple times with diff initial
conditions

global optimum

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

different local optima

Andrew Ng

## Random initialization

For i = 1 to 100 {  → Run K-means 100 times

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

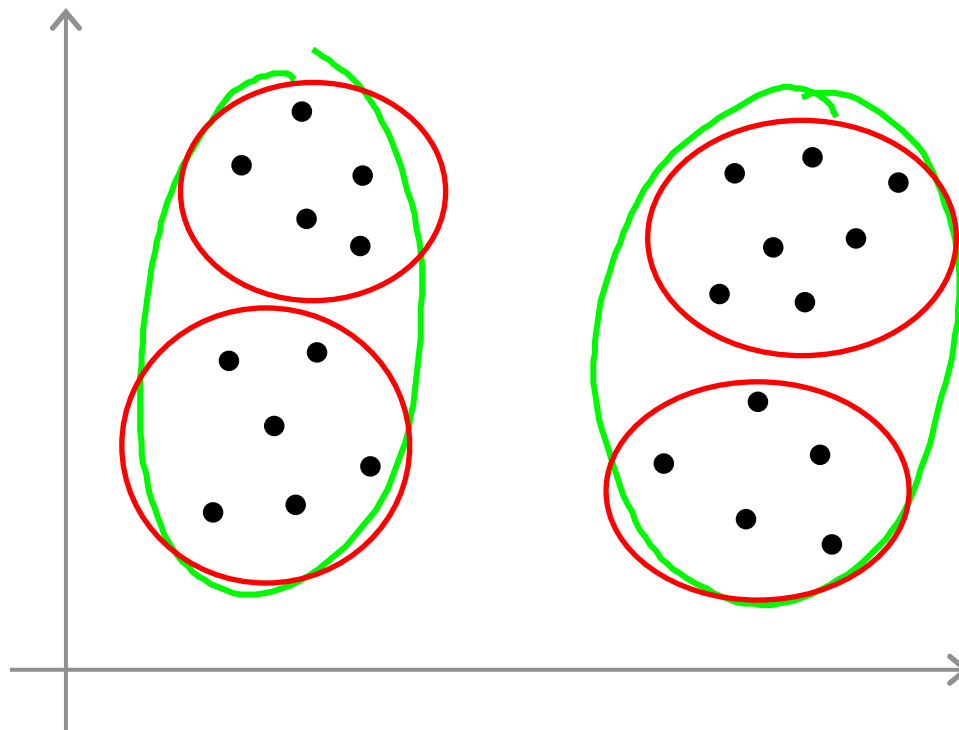Pick clustering that gave lowest cost $J(c^{(1)}, \ldots, c^{(m)}, \mu_1, \ldots, \mu_K)$

# Clustering

## Choosing the number of clusters

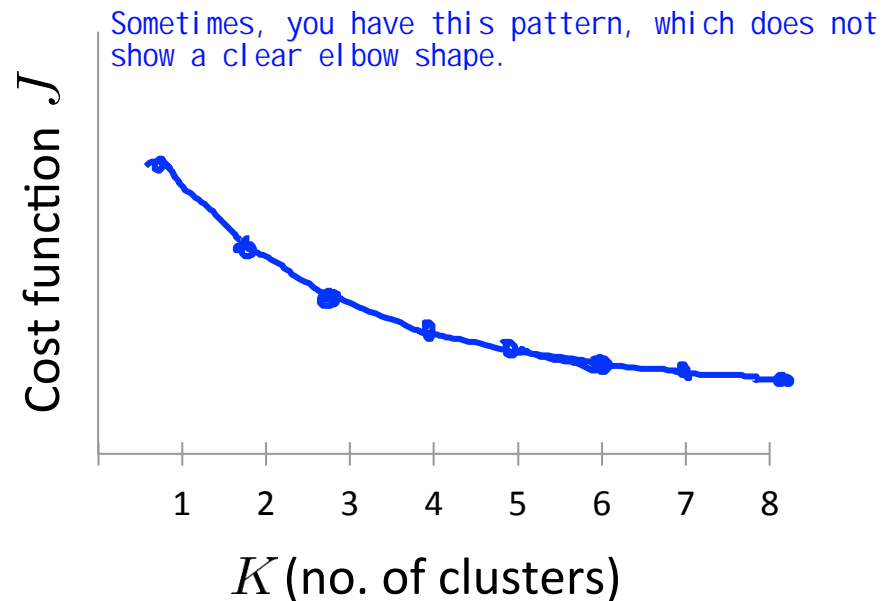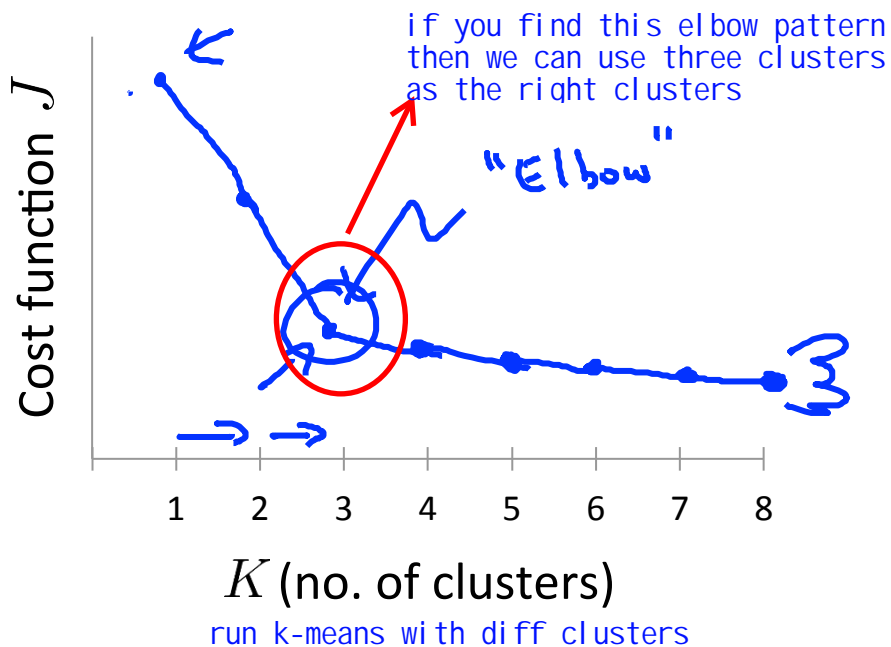Machine Learning

# What is the right value of K?

Choosing the value of K

Elbow method: → No high expectation for any particular problem

Cost function $J$

if you find this elbow pattern
then we can use three clusters
as the right clusters

"Elbow"

$K$ (no. of clusters)

run k-means with diff clusters

Cost function $J$

Sometimes, you have this pattern, which does not
show a clear elbow shape.

$K$ (no. of clusters)

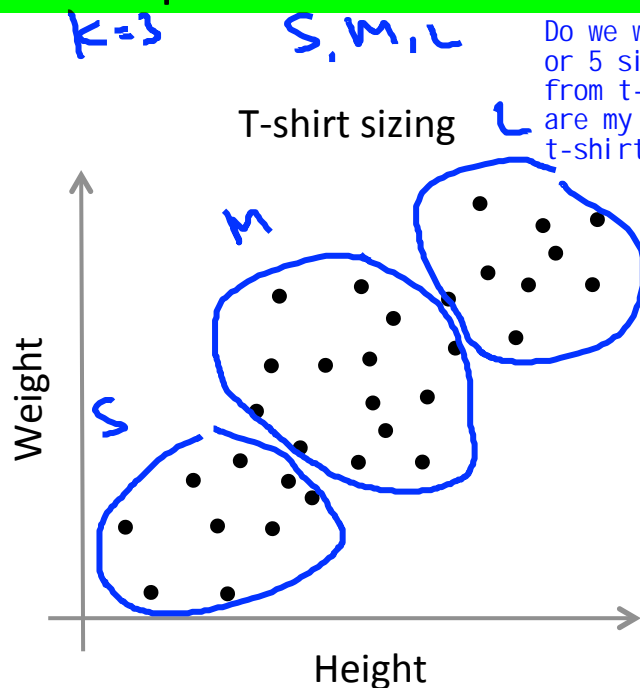Andrew Ng

**Choosing the value of K** <span style="color:blue">Another way to choose number of clusters</span>

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.

<span style="color:blue">K=3      S, M, L</span>

E.g.

T-shirt sizing

<span style="color:blue">Do we want to choose 3 sizes or 5 sizes? We can decide this from t-shirt sale business. e.g. are my customer happier with 5 t-shirt sizes?</span>

<span style="color:blue">K=5      XS, S, M, L, XL</span>

T-shirt sizing