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

# 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

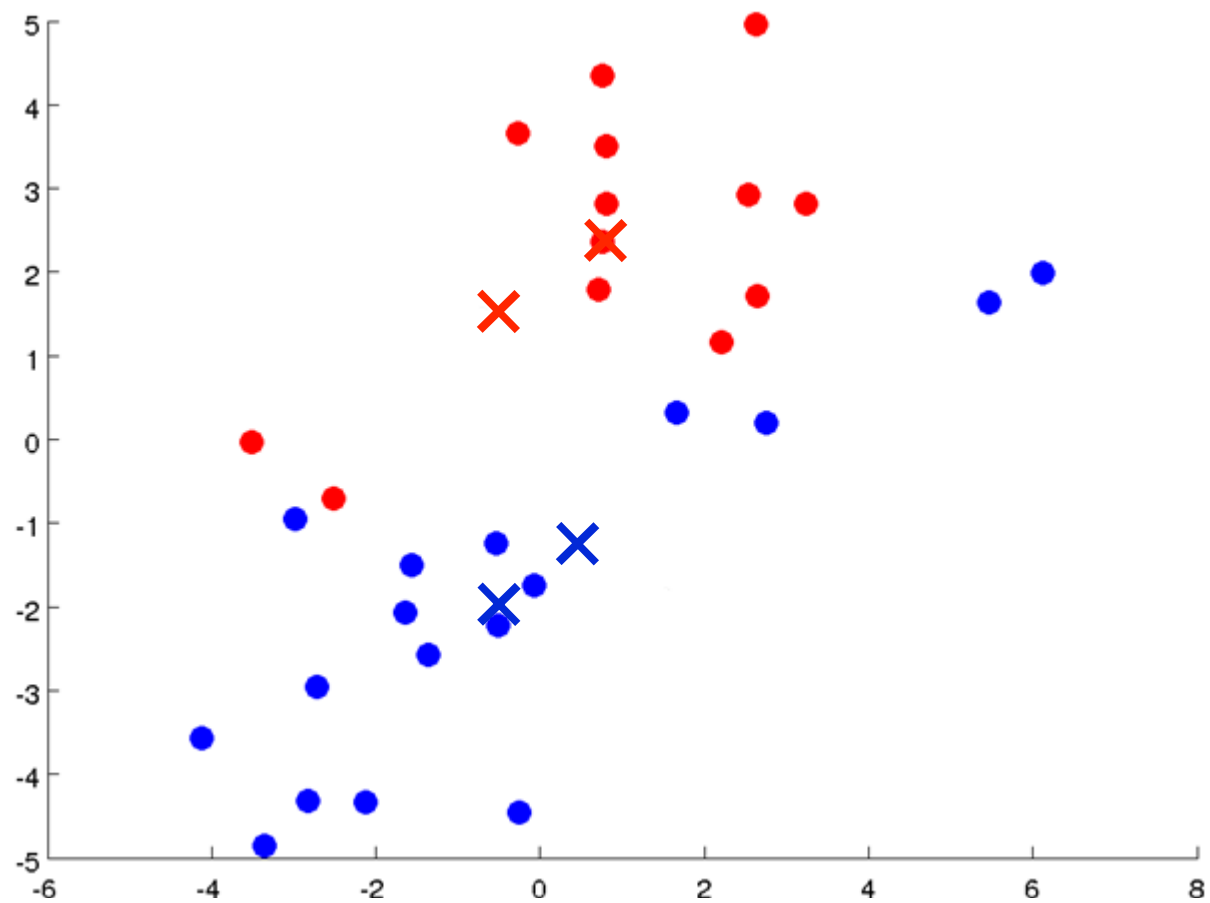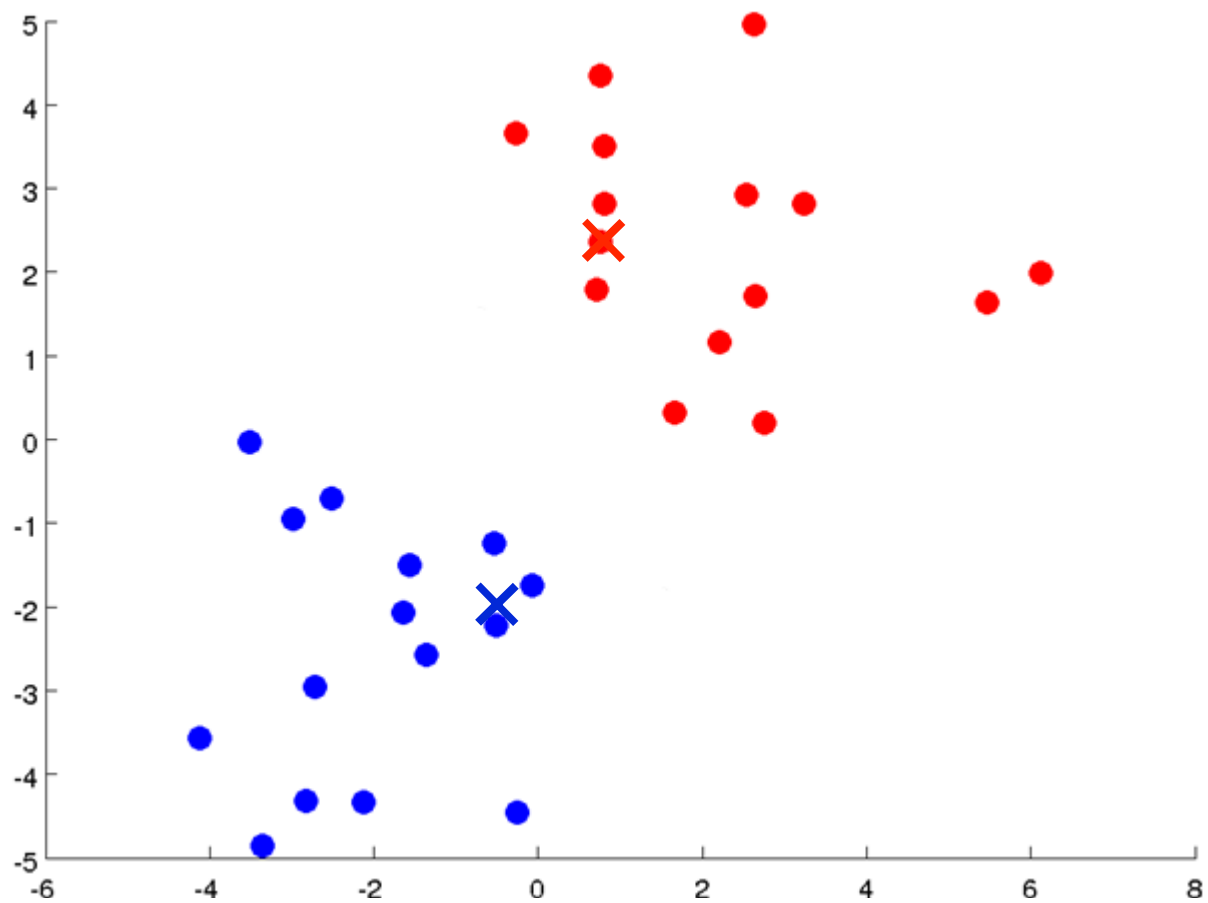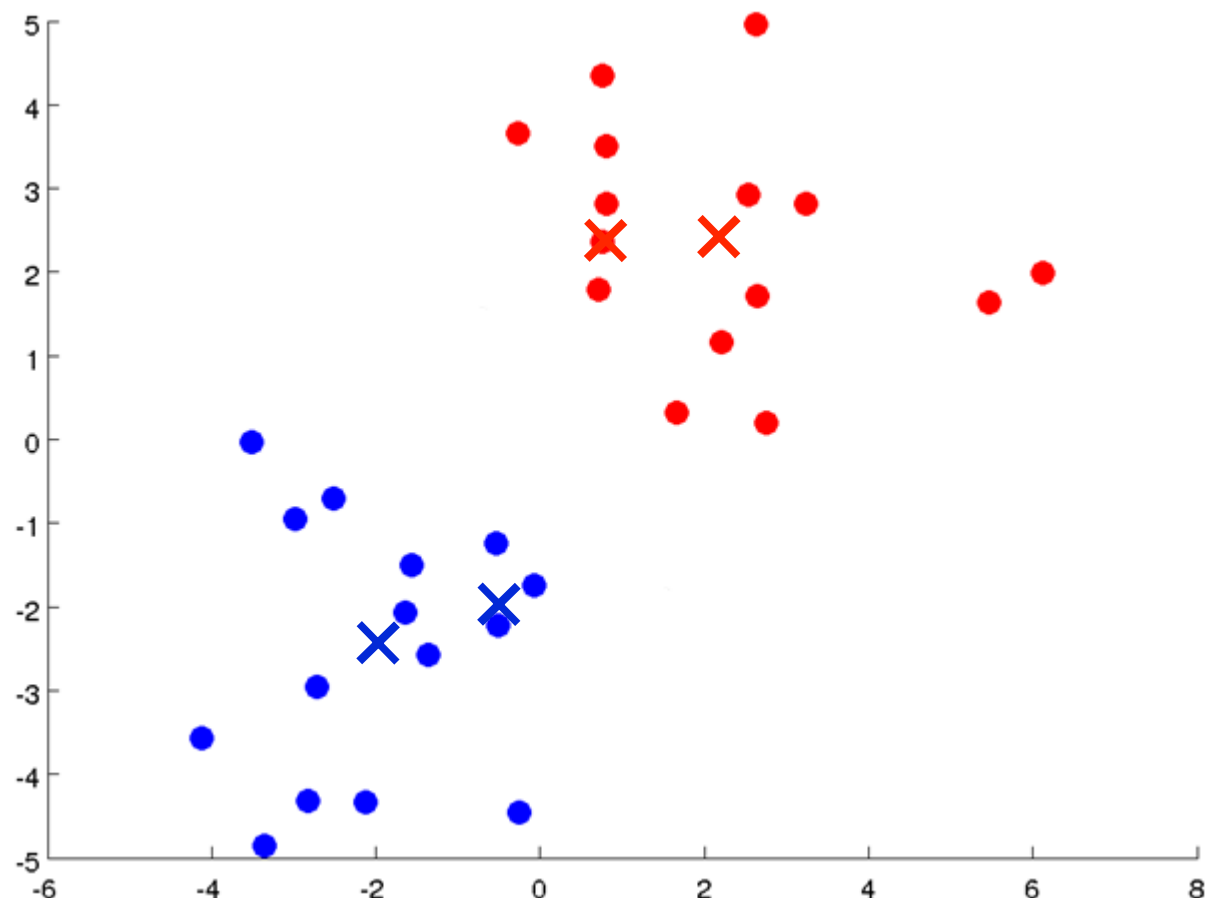cluster centroids

Andrew Ng

# K-means algorithm

Input:

- $K$ (number of clusters)

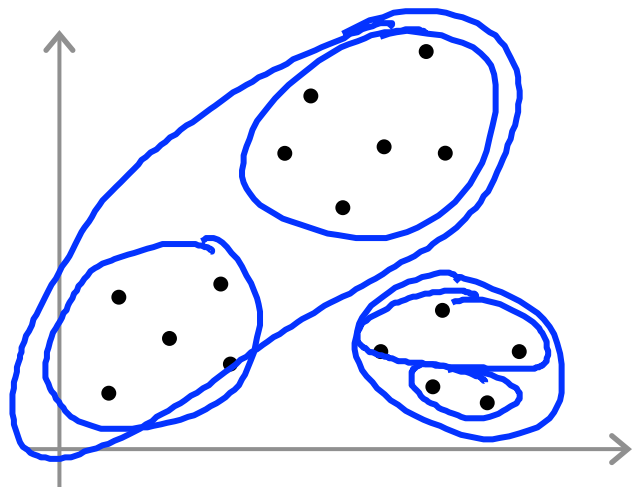- Training set $\{x^{(1)}, x^{(2)}, \ldots, x^{(m)}\}$

$x^{(i)} \in \mathbb{R}^n$ (drop $x_0 = 1$ convention)
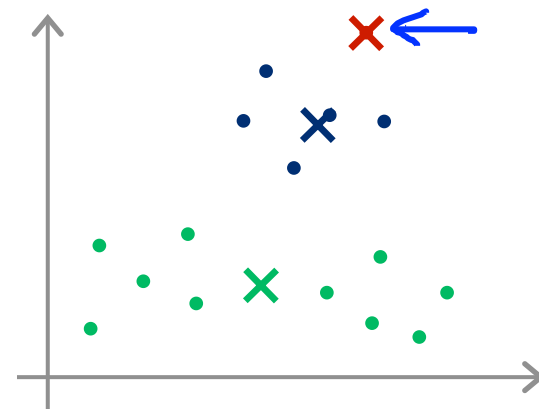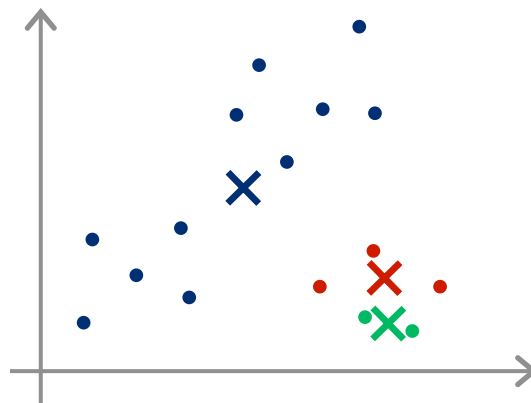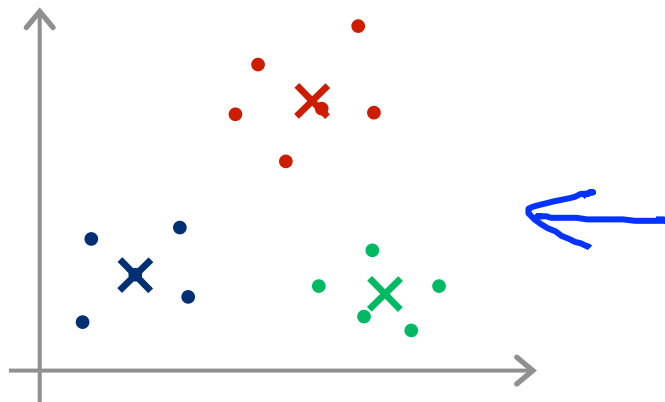
**K-means algorithm**

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

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 $k$
      }

# Local optima



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

Andrew Ng

**Random initialization**

For i = 1 to 100 {

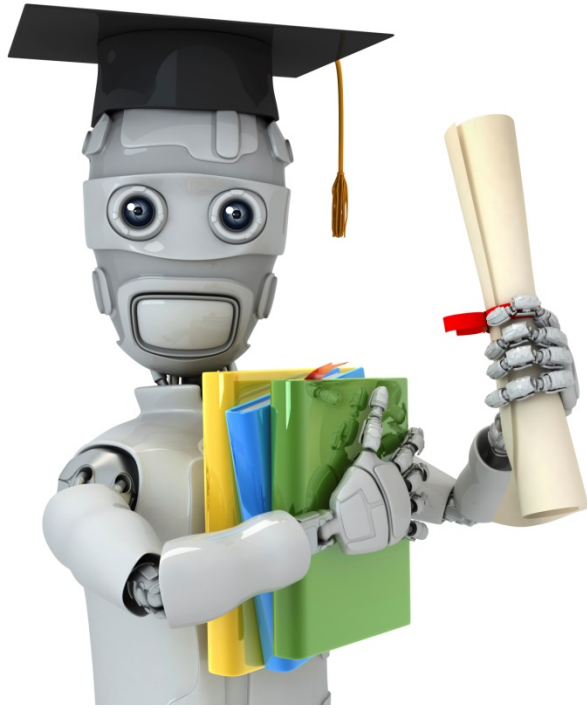       Randomly initialize K-means.
       Run K-means. Get $c^{(1)}, \ldots, c^{(m)}, \mu_1, \ldots, \mu_K$.
       Compute cost function (distortion)
         $J(c^{(1)}, \ldots, c^{(m)}, \mu_1, \ldots, \mu_K)$
       }


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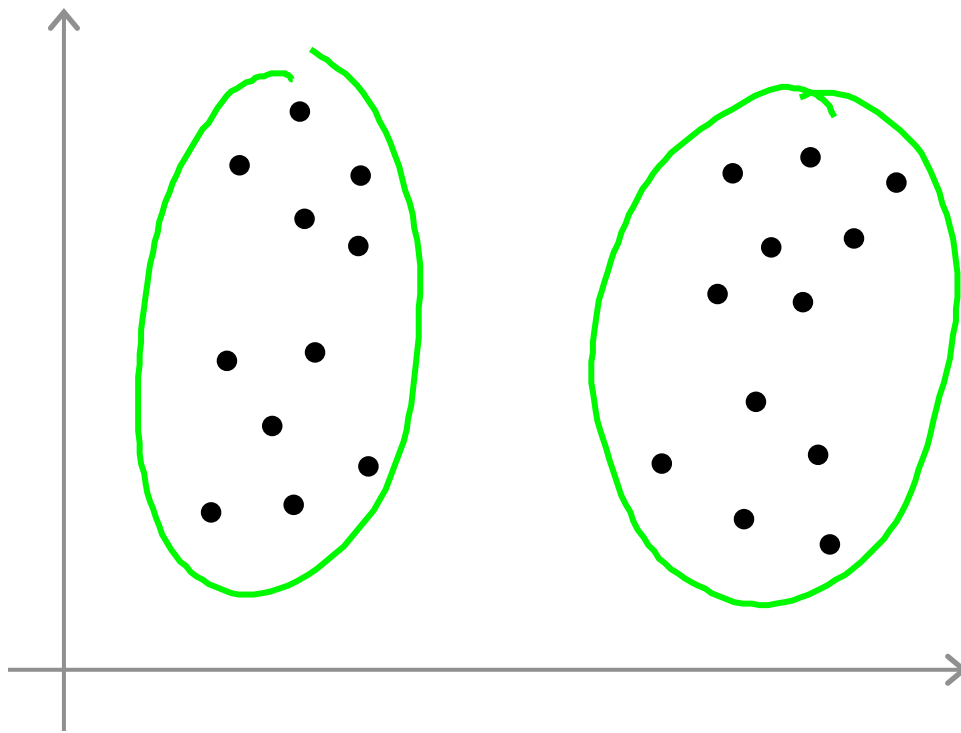
# Clustering

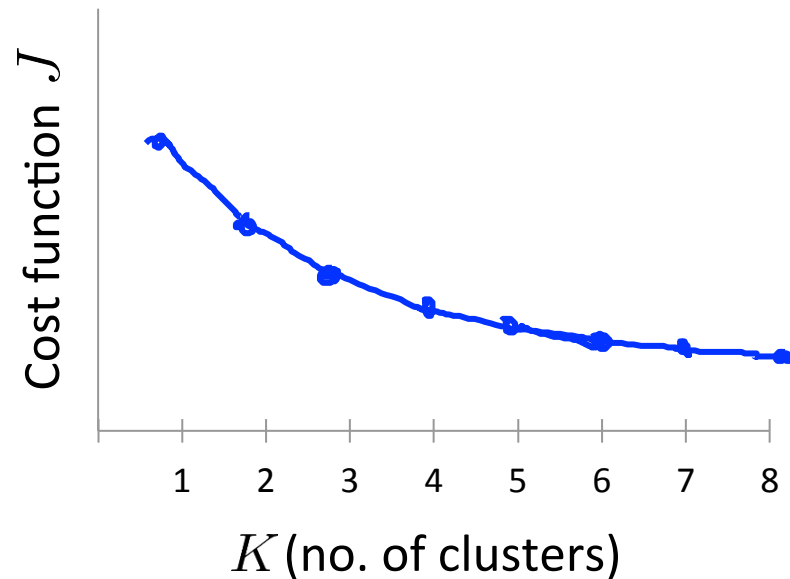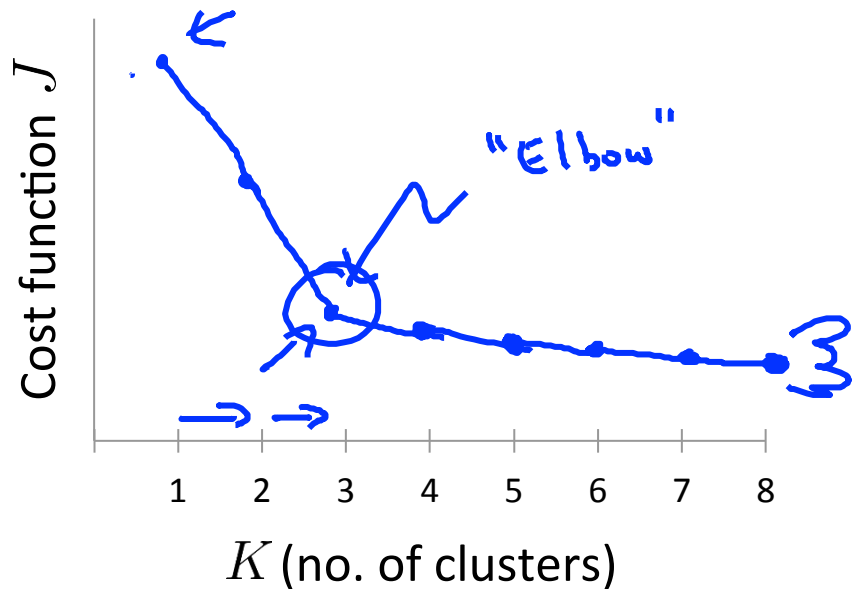## Choosing the number of clusters

Machine Learning

# What is the right value of K?

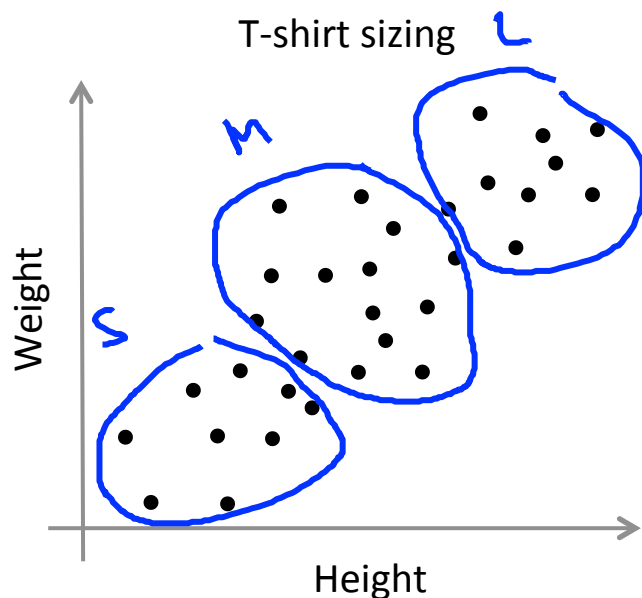# Choosing the value of K

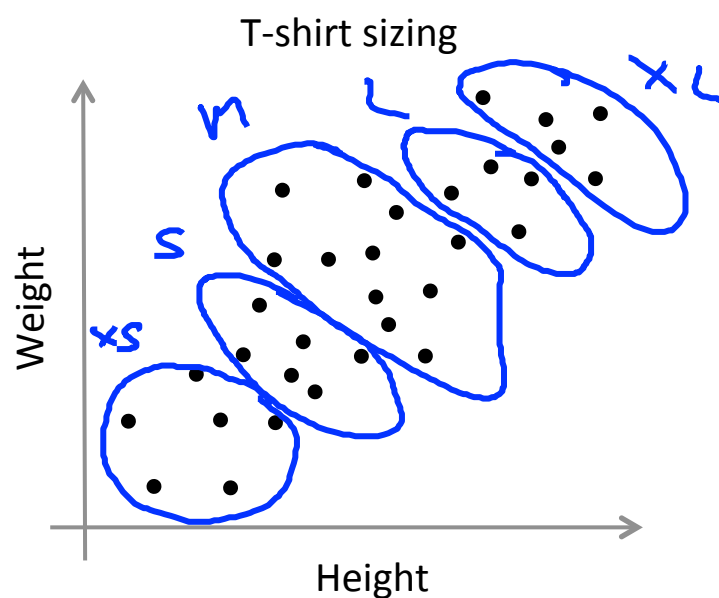Elbow method:

# 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

Weight

Height



T-shirt sizing

Weight

Height

# DATA MINING CLUSTERING

*by Panayiotis Tsaparas*

The k-means algorithm

Hierarchical Clustering

The DBSCAN algorithm
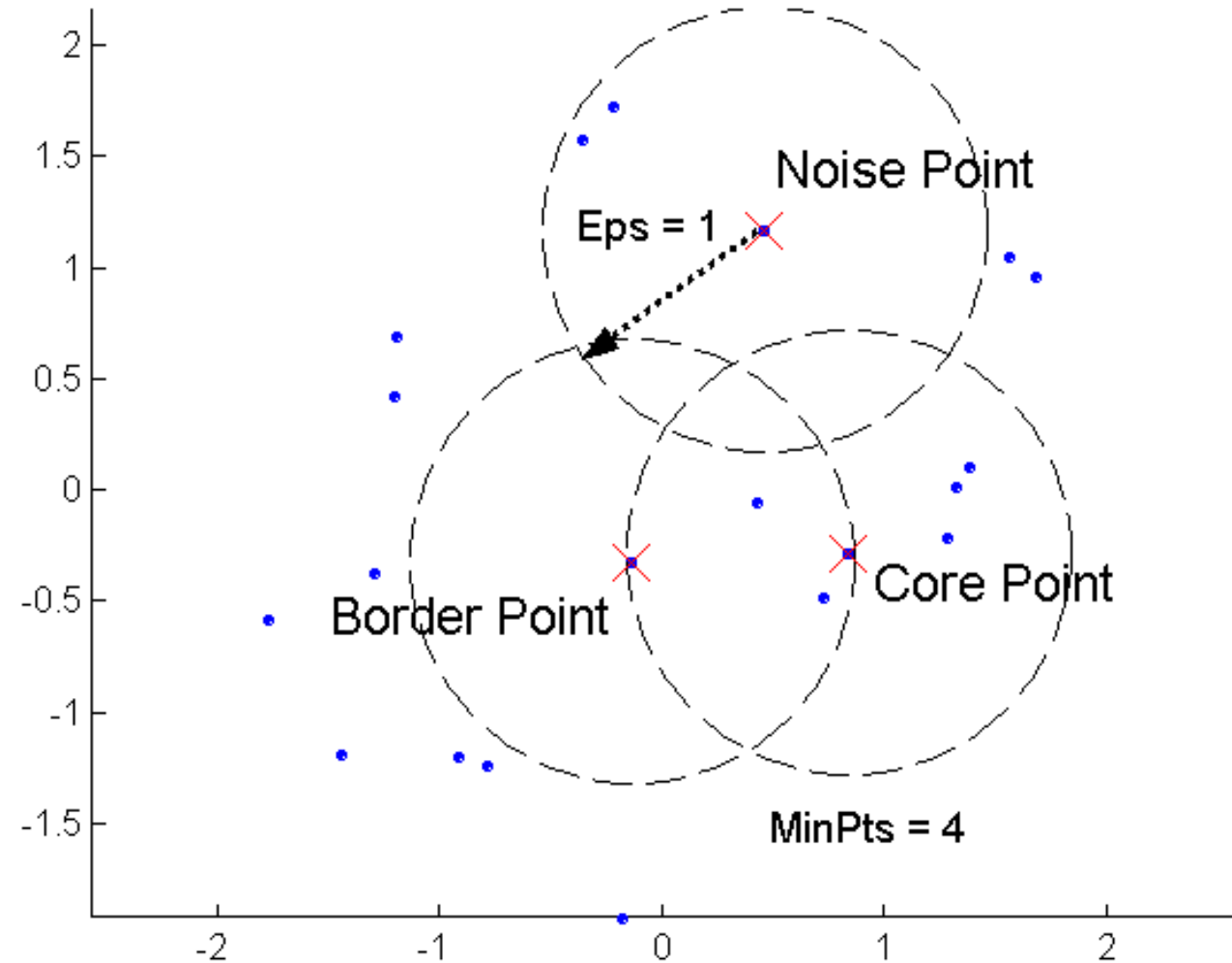
Evaluation

# DBSCAN

# DBSCAN: Density-Based Clustering

- DBSCAN is a Density-Based Clustering algorithm

- Reminder: In density-based clustering we partition points into dense regions separated by not-so-dense regions.

- Important Questions:
  - How do we measure density?
  - What is a dense region?

- DBSCAN:
  - Density at point p: number of points within a circle of radius Eps
  - Dense Region: A circle of radius Eps that contains at least MinPts points

# DBSCAN

- Characterization of points
  - A point is a core point if it has more than a specified number of points (MinPts) within Eps
    - These points belong in a dense region and are at the interior of a cluster

  - A border point has fewer than MinPts within Eps, but is in the neighborhood of a core point.

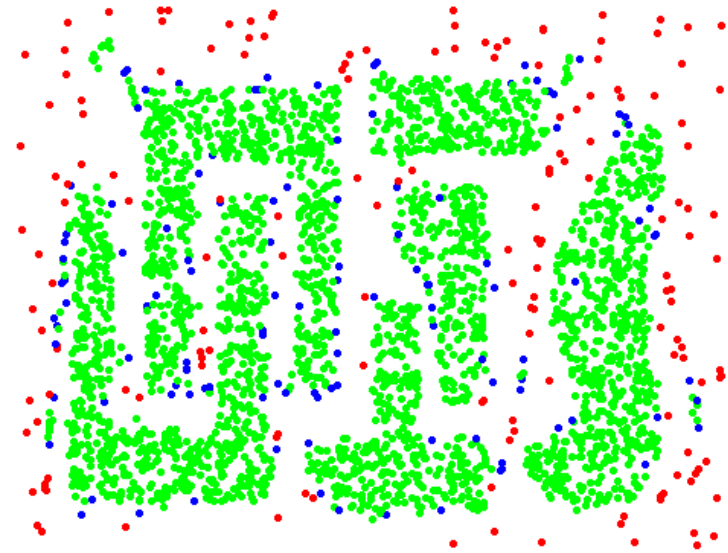  - A noise point is any point that is not a core point or a border point.

# DBSCAN: Core, Border, and Noise Points

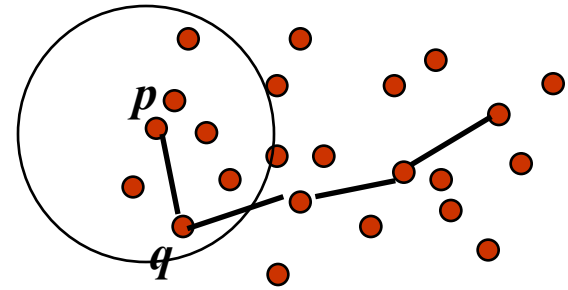# DBSCAN: Core, Border and Noise Points



Original Points

Point types: core, border and noise

Eps = 10, MinPts = 4

# Density-Connected points

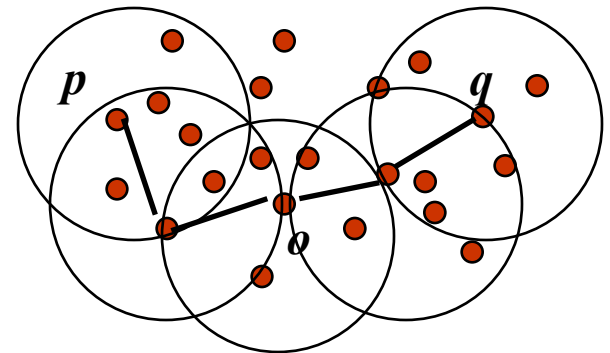- Density edge

  - We place an edge between two core points q and p if they are within distance Eps.



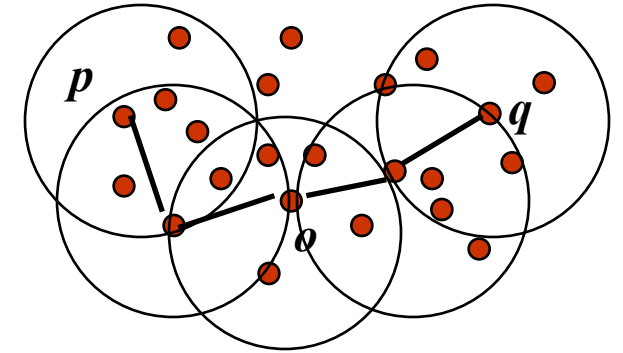- Density-connected

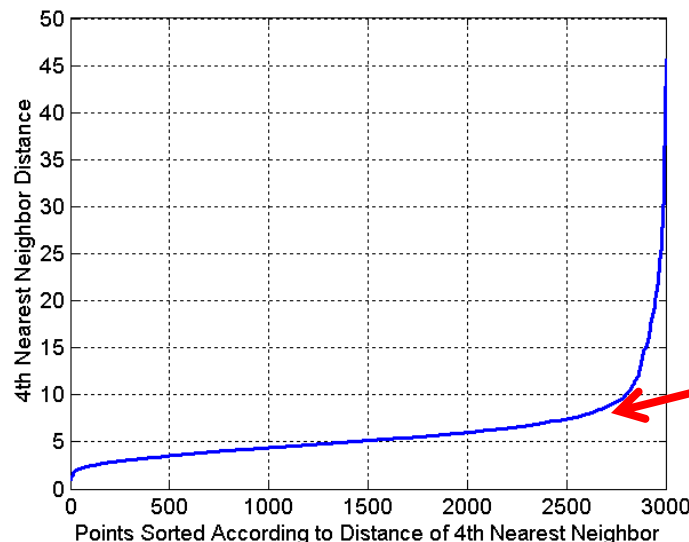  - A point p is density-connected to a point q if there is a path of edges from p to q

# DBSCAN Algorithm



- Label points as core, border and noise

- Eliminate noise points

- For every core point q that has not been assigned to a cluster

  - Create a new cluster with the point q and all the points that are density-connected to q.

- Assign border points to the cluster of the closest core point.

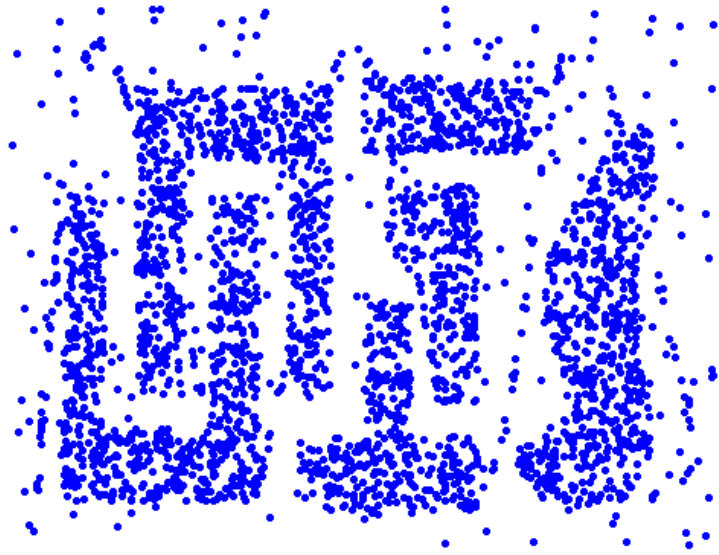# DBSCAN: Determining Eps and MinPts

- Try different minPts= k
- So, plot sorted distance of every point to its $k^{th}$ nearest neighbor
- Find the distance d where there is a "knee" in the curve
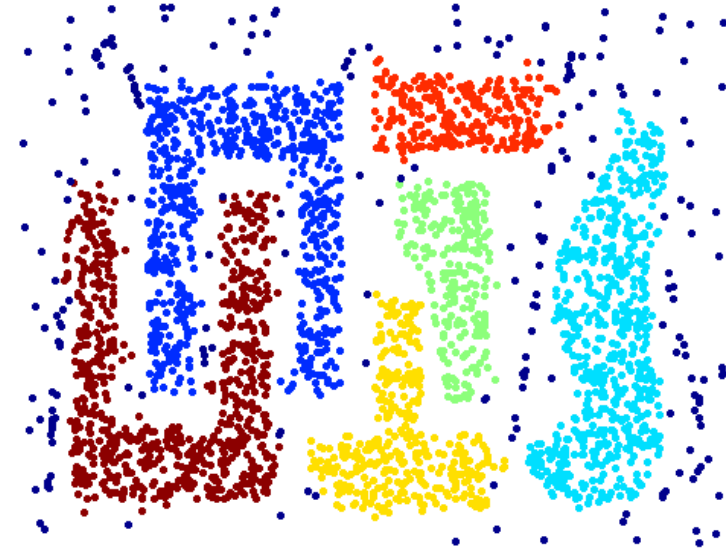  - Eps = d, MinPts = k



Noise points have the $k^{th}$ nearest neighbor at farther distance

Eps ~ 7-10
MinPts = 4

# When DBSCAN Works Well



Original Points

Clusters

- Resistant to Noise

- Can handle clusters of different shapes and sizes

# DBSCAN: Sensitive to Parameters



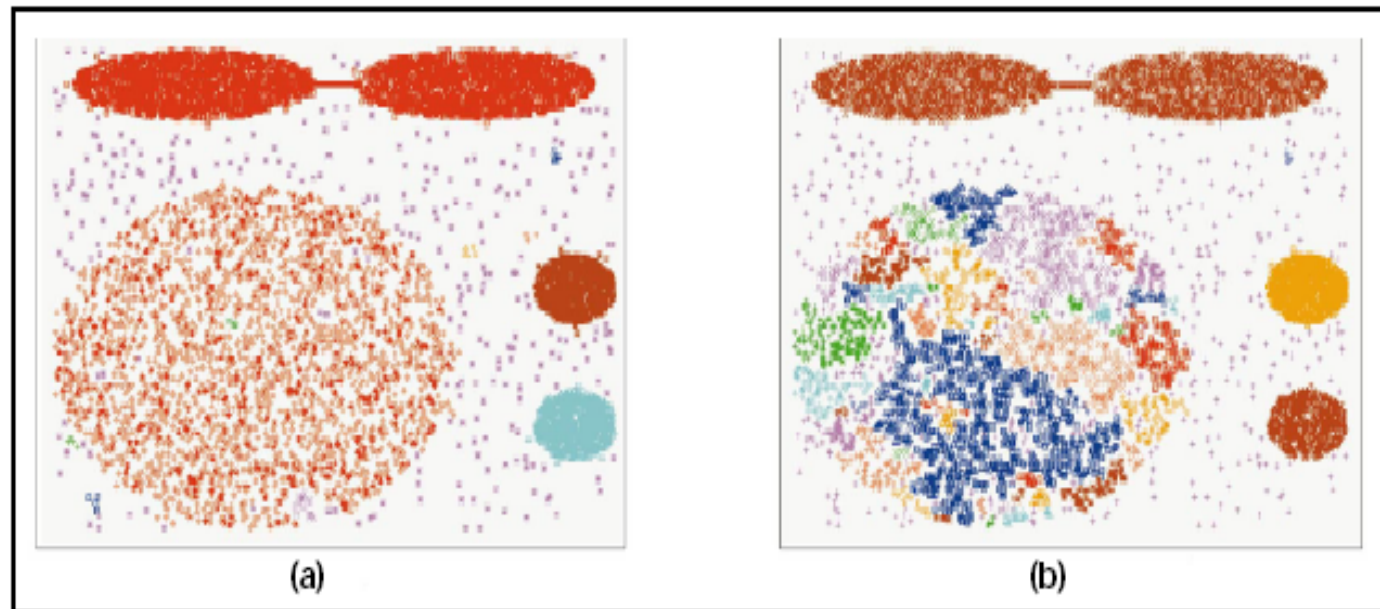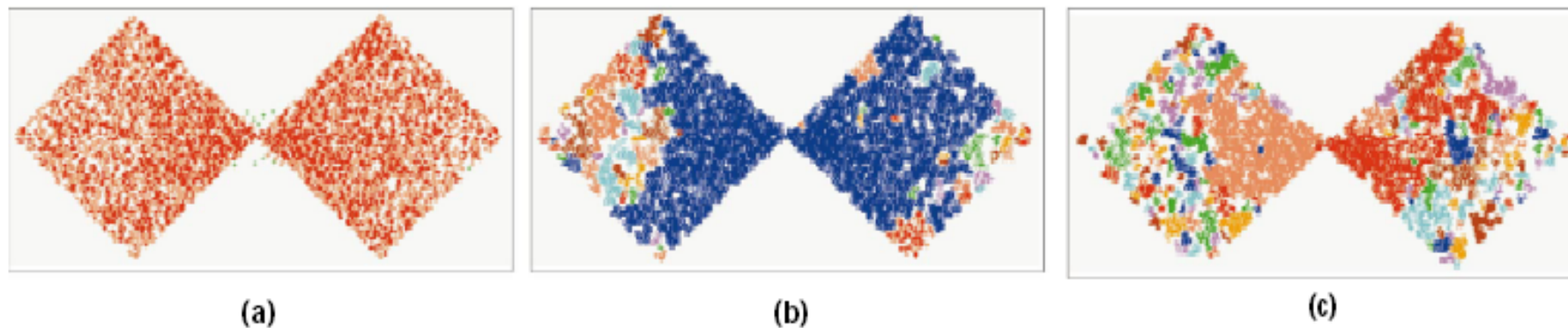Figure 8. DBScan results for DS1 with MinPts at 4 and Eps at (a) 0.5 and (b) 0.4.
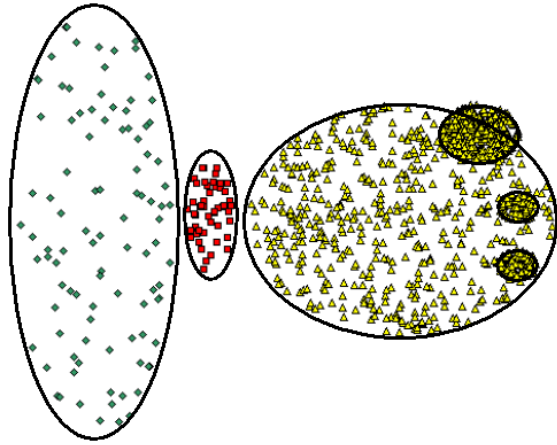
Figure 9. DBScan results for DS2 with MinPts at 4 and Eps at (a) 5.0, (b) 3.5, and (c) 3.0.

# When DBSCAN Does NOT Work Well



Original Points



(MinPts=4, Eps=9.75).



(MinPts=4, Eps=9.92)

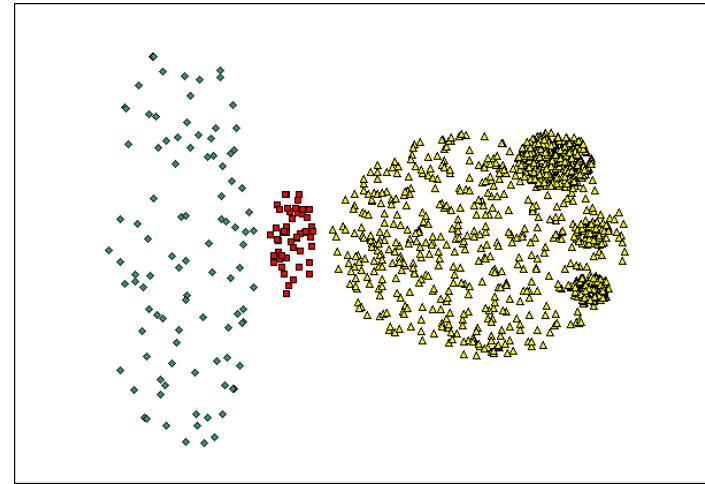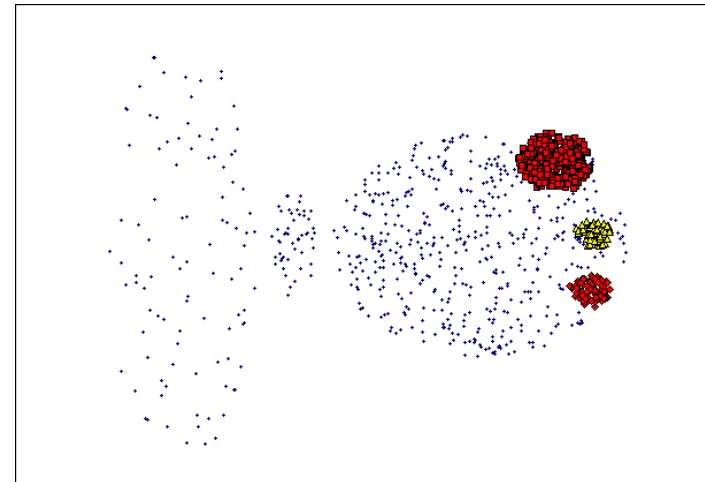- ## Varying densities

  - large differences in densities, since the minPts-Eps combination cannot then be chosen appropriately for all clusters.

- ## High-dimensional data

  - difficult to find an appropriate value for Eps