# CSE 4621
# Machine Learning

Lecture 14

**Md. Hasanul Kabir, PhD.**

Professor, CSE Department

Islamic University of Technology (IUT)

# What is Cluster Analysis?

- Cluster: A collection of data objects
  - similar (or related) to one another within the same group
  - dissimilar (or unrelated) to the objects in other groups
- Cluster analysis (or *clustering*, *data segmentation, …*)
  - Finding similarities between data according to the characteristics found in the data and grouping similar data objects into clusters
- Unsupervised learning: no predefined classes (i.e., *learning by observations* vs. learning by examples: supervised)
- Typical applications
  - As a stand-alone tool to get insight into data distribution
  - As a preprocessing step for other algorithms

# Clustering for Data Understanding and Applications

- Biology: taxonomy of living things: kingdom, phylum, class, order, family, genus and species

- Information retrieval: document clustering

- Land use: Identification of areas of similar land use in an earth observation database

- Marketing: Help marketers discover distinct groups in their customer bases, and then use this knowledge to develop targeted marketing programs

- City-planning: Identifying groups of houses according to their house type, value, and geographical location

- Earth-quake studies: Observed earth quake epicenters should be clustered along continent faults

- Web Search: Clustering can be used to organize the search results into groups and present the results in a concise and easily accessible way.

- Information Retrieval: Cluster documents into topics.

# Clustering as a Preprocessing Tool (Utility)

- Summarization:
  - Preprocessing for regression, PCA, classification, and association analysis
- Compression:
  - Image processing: vector quantization
- Finding K-nearest Neighbors
  - Localizing search to one or a small number of clusters
- Outlier detection
  - Outliers are often viewed as those "far away" from any cluster

# Quality: What Is Good Clustering?

- A good clustering method will produce high quality clusters

  - high intra-class similarity: cohesive within clusters

  - low inter-class similarity: distinctive between clusters

- The quality of a clustering method depends on

  - the similarity measure used by the method

  - its implementation, and

  - Its ability to discover some or all of the hidden patterns

# Measure the Quality of Clustering

- Dissimilarity/Similarity metric
  - Similarity is expressed in terms of a distance function, typically metric: $d(i, j)$
  - The definitions of distance functions are usually rather different for interval-scaled, boolean, categorical, ordinal ratio, and vector variables
  - Weights should be associated with different variables based on applications and data semantics
- Quality of clustering:
  - There is usually a separate "quality" function that measures the "goodness" of a cluster.
  - It is hard to define "similar enough" or "good enough"
    - The answer is typically highly subjective

# Considerations for Cluster Analysis

- Partitioning criteria
  - Single level vs. hierarchical partitioning (often, multi-level hierarchical partitioning is desirable). E.g. Politics, Sports: Football, Cricket, volleyball, etc.

- Separation of clusters
  - Exclusive (e.g., one customer belongs to only one region) vs. non-exclusive (e.g., one document may belong to more than one class)

- Similarity measure
  - Distance-based (e.g., Euclidian, road network, vector)  vs. connectivity-based (e.g., density or contiguity)
  - Distance-based methods can often take advantage of optimization techniques, density- and continuity-based methods can often find clusters of arbitrary shape

- Clustering space
  - Full space (often when low dimensional) vs. subspaces (often in high-dimensional clustering)

# Major Clustering Approaches

- Partitioning approach:
  - Construct various partitions and then evaluate them by some criterion, e.g., minimizing the sum of square errors
  - Typical methods: k-means, k-medoids, CLARANS
- Hierarchical approach:
  - Create a hierarchical decomposition of the set of data (or objects) using some criterion
  - Typical methods: Diana, Agnes, BIRCH, CAMELEON
- Density-based approach:
  - Based on connectivity and density functions
  - Typical methods: DBSACN, OPTICS, DenClue
- Grid-based approach:
  - based on a multiple-level granularity structure
  - quantize object space into a finite number of cells (grid structure)
  - Typical methods: STING, WaveCluster, CLIQUE

# Overview of Clustering Methods

| Method | General Characteristics |
|---|---|
| Partitioning methods | – Find mutually exclusive clusters of spherical shape<br>– Distance-based<br>– May use mean or medoid (etc.) to represent cluster center<br>– Effective for small- to medium-size data sets |
| Hierarchical methods | – Clustering is a hierarchical decomposition (i.e., multiple levels)<br>– Cannot correct erroneous merges or splits<br>– May incorporate other techniques like microclustering or consider object "linkages" |
| Density-based methods | – Can find arbitrarily shaped clusters<br>– Clusters are dense regions of objects in space that are separated by low-density regions<br>– Cluster density: Each point must have a minimum number of points within its "neighborhood"<br>– May filter out outliers |
| Grid-based methods | – Use a multiresolution grid data structure<br>– Fast processing time (typically independent of the number of data objects, yet dependent on grid size) |

# Partitioning Algorithms: Basic Concept

- Partitioning method: Partitioning a database **D** of **n** objects into a set of **k** clusters, such that the sum of squared distances is minimized (where $c_i$ is the centroid or medoid of cluster $C_i$)

$$E = \Sigma_{i=1}^{k} \Sigma_{p \in C_i} \| p - c_i \|^2$$

- Given *k <=n*, find a partition of *k clusters* that optimizes the chosen partitioning criterion
  - Global optimal: exhaustively enumerate all partitions
  - Heuristic methods: *k-means* and *k-medoids* algorithms
  - *k-means* (MacQueen'67, Lloyd'57/'82): Each cluster is represented by the center of the cluster
  - *k-medoids* or PAM (Partition around medoids) (Kaufman & Rousseeuw'87): Each cluster is represented by one of the objects in the cluster

# The *K-Means* Clustering Method

- Given *k*, the *k-means* algorithm is implemented as

**Algorithm: k-means.** The *k*-means algorithm for partitioning, where each cluster's center is represented by the mean value of the objects in the cluster.

**Input:**

- *k*: the number of clusters,
- *D*: a data set containing *n* objects.

**Output:** A set of *k* clusters.

**Method:**

(1) arbitrarily choose *k* objects from *D* as the initial cluster centers;
(2) **repeat**
(3)     (re)assign each object to the cluster to which the object is the most similar, based on the mean value of the objects in the cluster;
(4)     update the cluster means, that is, calculate the mean value of the objects for each cluster;
(5) **until** no change;

# An Example of *K-Means* Clustering

# Clustering

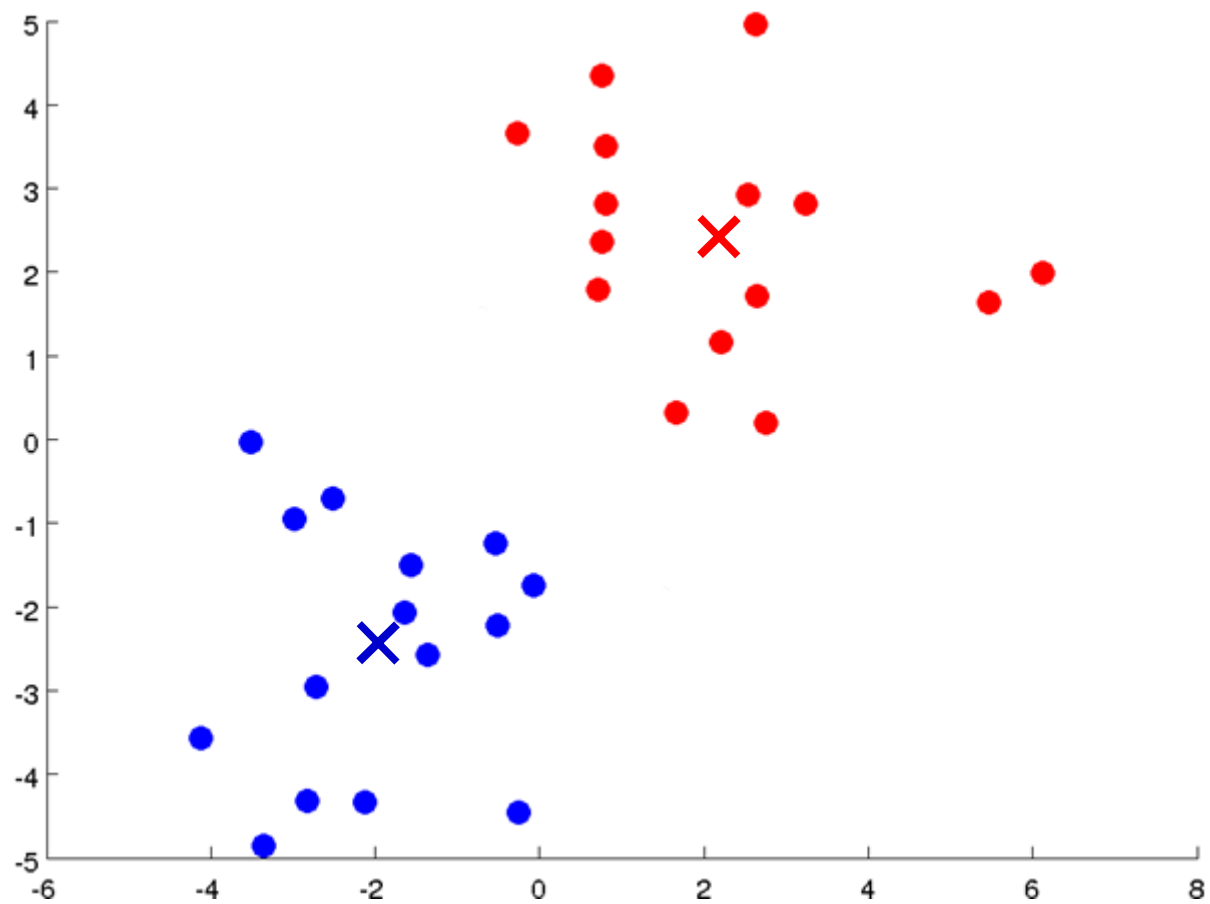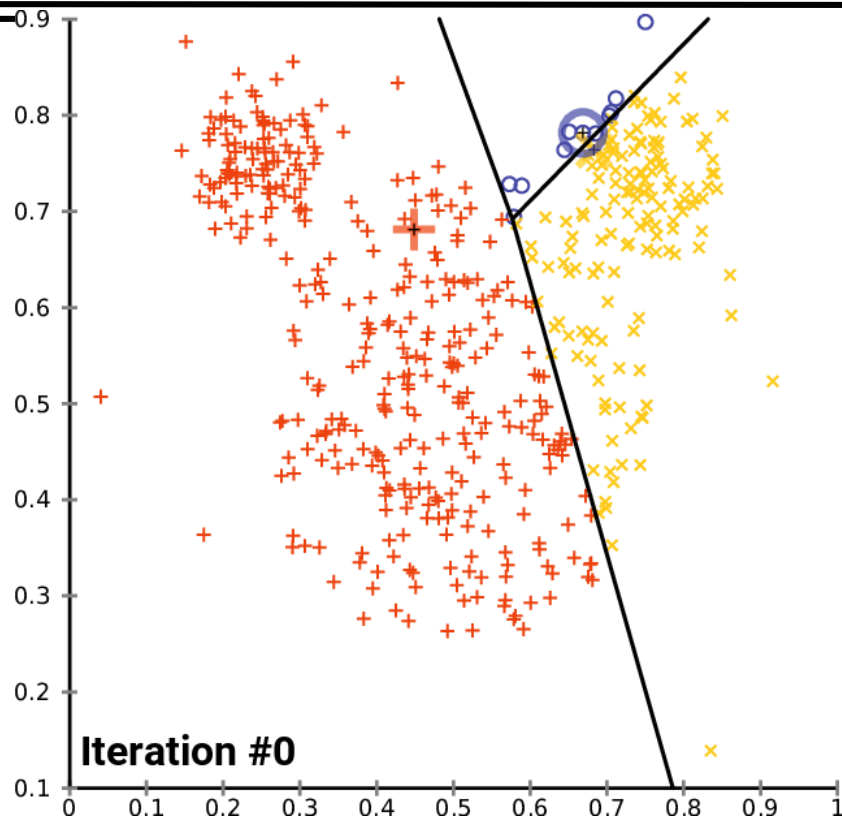## K-means Example

Machine Learning

cluster centroids

Andrew Ng

Andrew Ng

# Graphical Example of *K-Means* Clustering



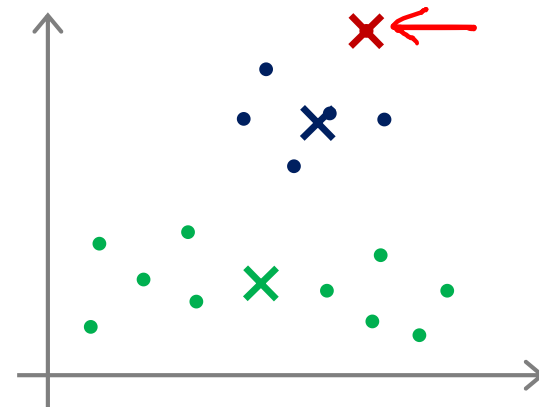Iteration #0
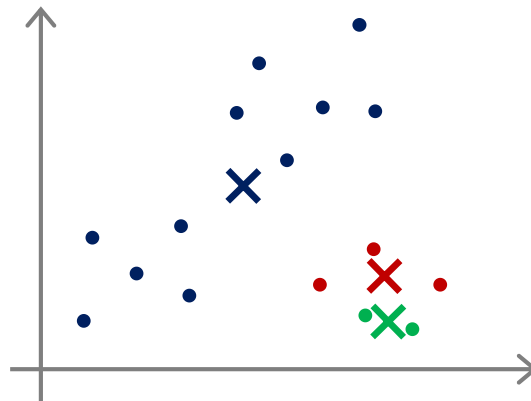
# Local optima

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

Andrew Ng

# Random initialization

For i = 1 to 100 {

      Randomly initialize K-means.
      Run K-means. Get $(c_1, c_2, ..., c_K)$ .
      Compute cost function (distortion)

$$E(c_1, c_2, ..., c_K)$$

      }

Pick clustering that gave lowest cost $E(c_1, c_2, ..., c_K)$

# Comments on the *K-Means* Method

- <u>Strength:</u> *Efficient*: $O(tkn)$, where $n$ is # objects, $k$ is # clusters, and $t$ is # iterations. Normally, $k$, $t$ << $n$.
    - Comparing: PAM: $O(k(n-k)^2)$, CLARA: $O(ks^2 + k(n-k))$
- <u>Comment:</u> Often terminates at a *local optimal*.
- <u>Weakness</u>
    - Applicable only to objects in a continuous n-dimensional space
        - Using the k-modes method for categorical data
        - In comparison, k-medoids can be applied to a wide range of data
    - Need to specify $k$, the *number* of clusters, in advance (there are ways to automatically determine the best k (see Hastie et al., 2009)
    - Sensitive to **noisy data and outliers**
    - Not suitable to discover clusters with *non-convex shapes*

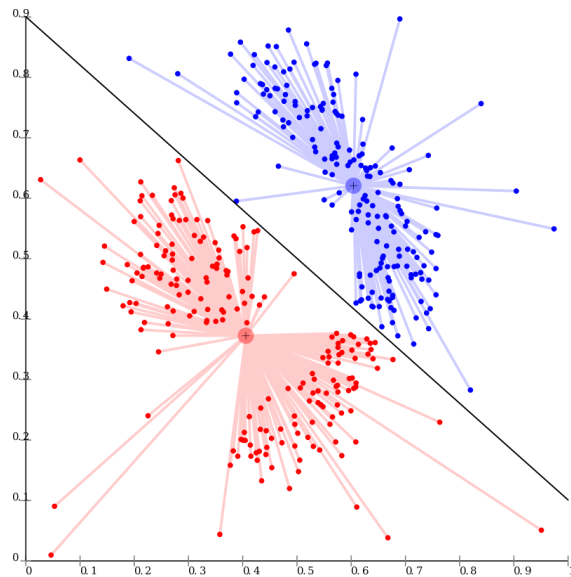# *k*-means cannot represent density-based clusters
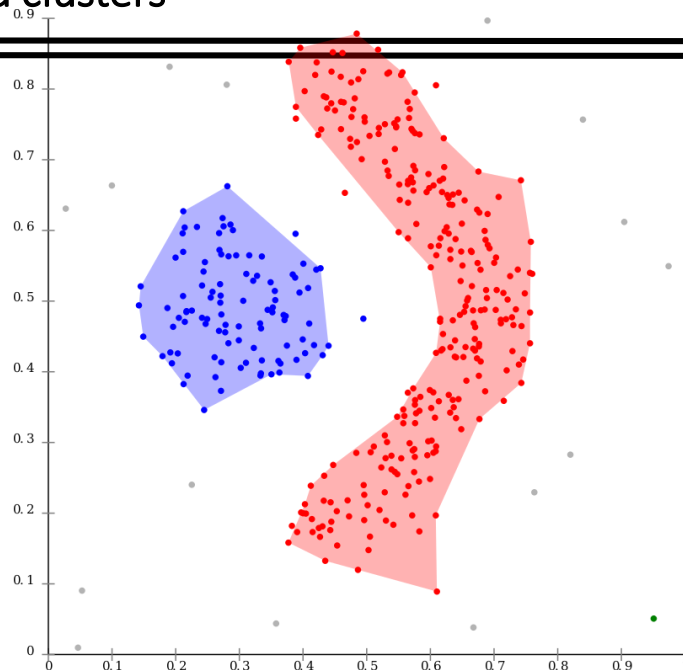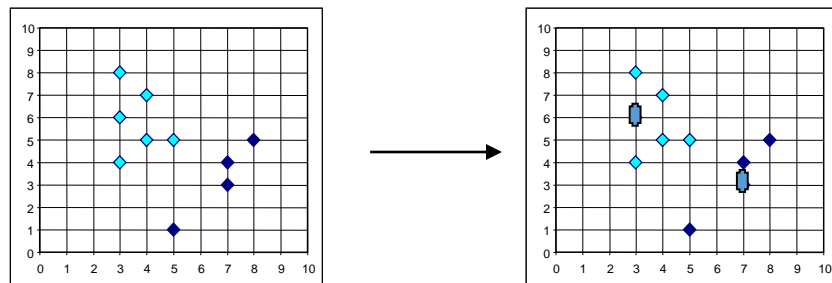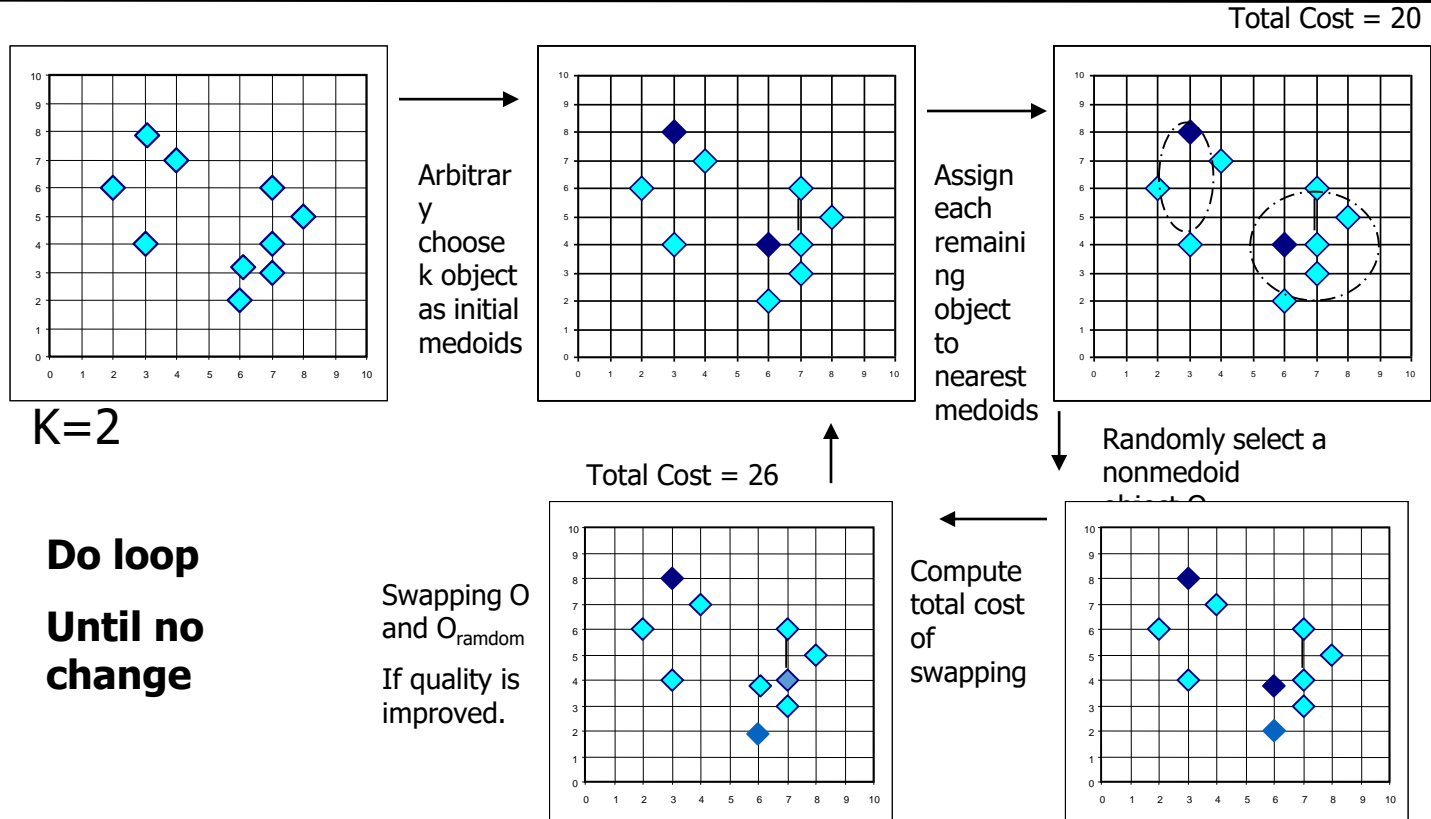


Fig: Convex Shaped

Fig: Non-Convex Shaped

# What Is the Problem of the K-Means Method?

- The k-means algorithm is sensitive to outliers !

  - Since an object with an extremely large value may substantially distort the distribution of the data

- K-Medoids:  Instead of taking the **mean** value of the object in a cluster as a reference point, **medoids** can be used, which is the **most centrally located** object in a cluster

# PAM: A Typical K-Medoids Algorithm



Total Cost = 20

K=2

Arbitrary choose k object as initial medoids

Assign each remaining object to nearest medoids

Randomly select a nonmedoid object O

Do loop

Until no change

Total Cost = 26

Swapping O and O_ramdom

If quality is improved.

Compute total cost of swapping

# The *K-Medoids* Clustering Method

**Algorithm:** *k*-**medoids.** PAM, a *k*-medoids algorithm for partitioning based on medoid or central objects.
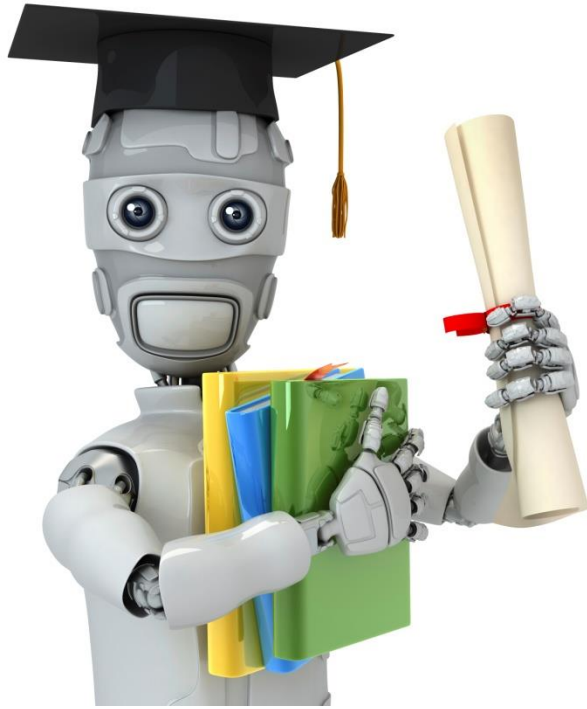
**Input:**

- *k*: the number of clusters,

- *D*: a data set containing *n* objects.

**Output:** A set of *k* clusters.

**Method:**

(1)  arbitrarily choose *k* objects in *D* as the initial representative objects or seeds;
(2)  **repeat**
(3)       assign each remaining object to the cluster with the nearest representative object;
(4)       randomly select a nonrepresentative object, $o_{random}$;
(5)       compute the total cost, *S*, of swapping representative object, $o_j$, with $o_{random}$;
(6)       **if** $S < 0$ **then** swap $o_j$ with $o_{random}$ to form the new set of *k* representative objects;
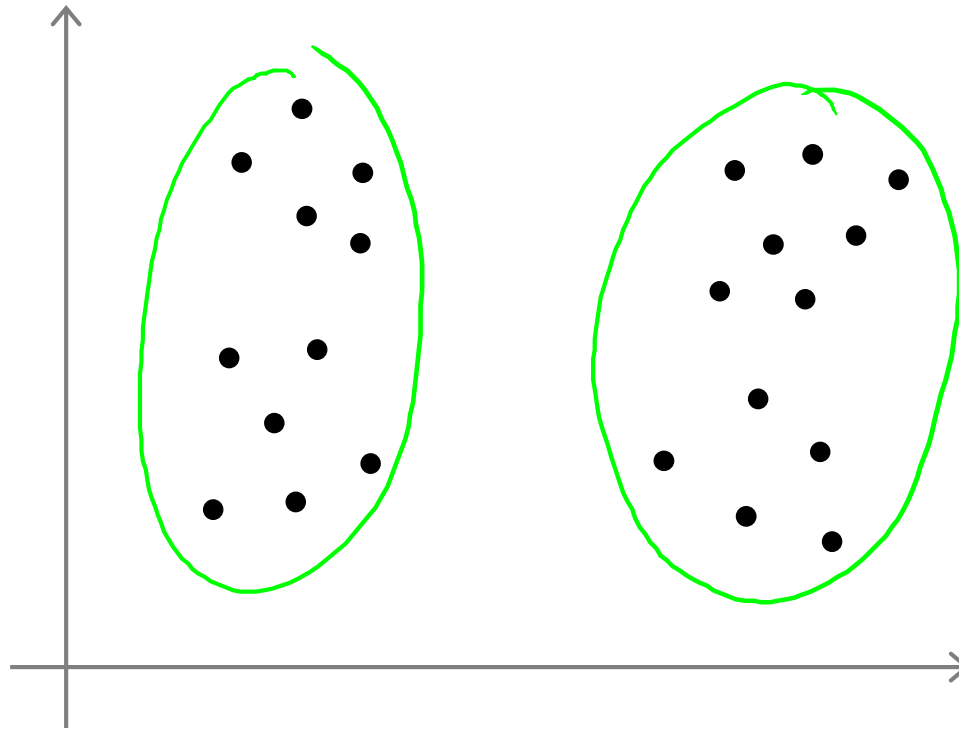(7)  **until** no change;

# Clustering
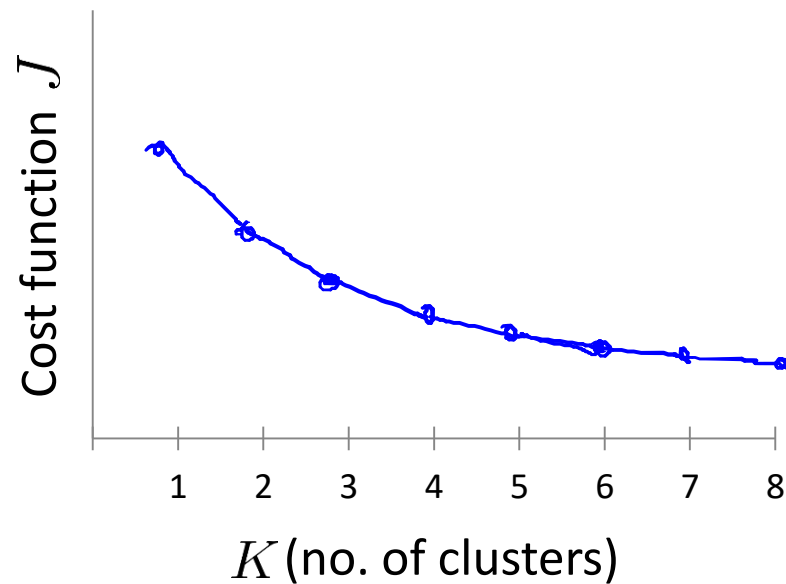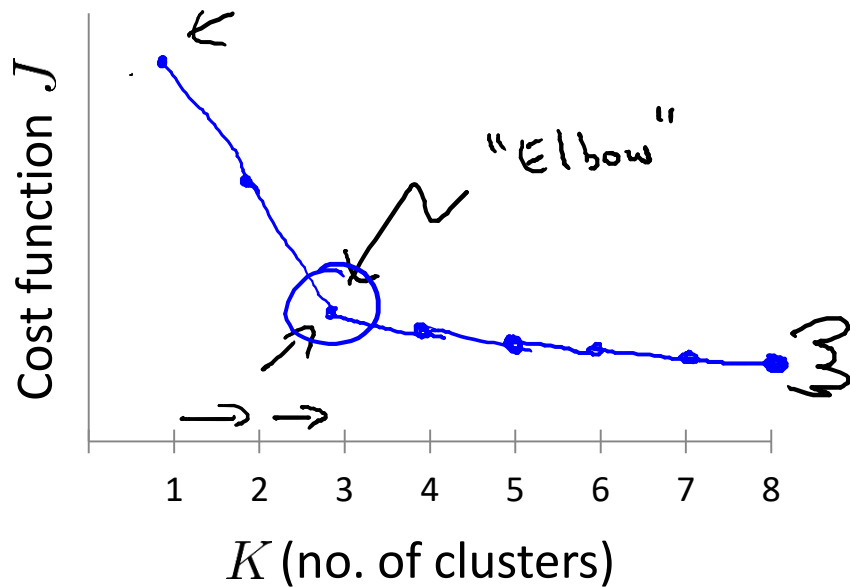
## Choosing the number of clusters

Machine Learning

# What is the right value of K?
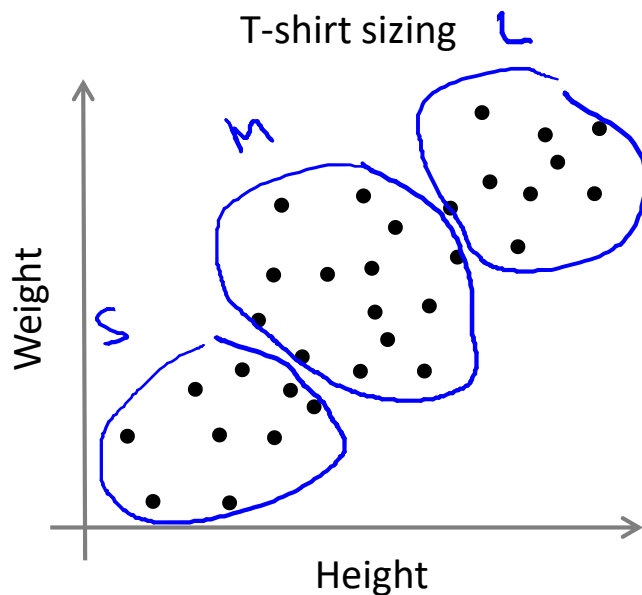
# Choosing the value of K

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