

# CLASSIFICATION AND CLUSTERING

Anveshi Charuvaka

# Learning from Data



- Classification
- Regression
- Clustering
- Anomaly Detection
- Contrast Set Mining

# Classification: Definition

- Given a collection of records (*training set*)
  - Each record contains a set of *attributes*, one of the attributes is the *class*.
- Find a *model* for class attribute as a function of the values of other attributes.
- Goal: previously unseen records should be assigned a class as accurately as possible.
  - A *test set* is used to determine the accuracy of the model. Usually, the given data set is divided into training and test sets, with training set used to build the model and test set used to validate it.

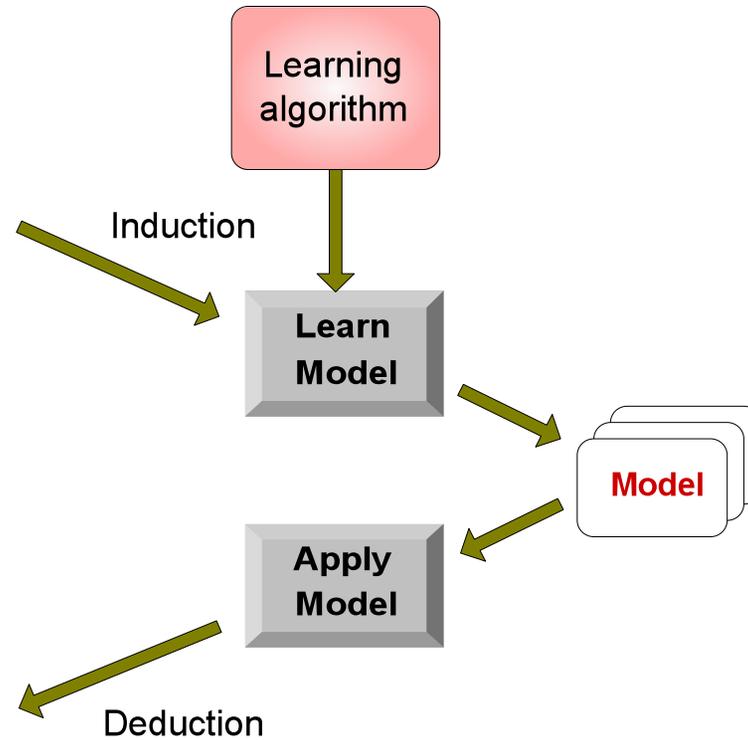
# Illustrating Classification Task

Tid	Attrib1	Attrib2	Attrib3	Class
1	Yes	Large	125K	No
2	No	Medium	100K	No
3	No	Small	70K	No
4	Yes	Medium	120K	No
5	No	Large	95K	Yes
6	No	Medium	60K	No
7	Yes	Large	220K	No
8	No	Small	85K	Yes
9	No	Medium	75K	No
10	No	Small	90K	Yes

Training Set

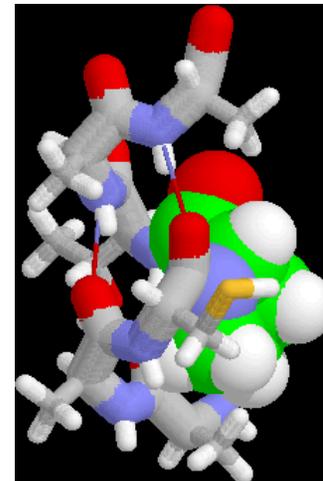
Tid	Attrib1	Attrib2	Attrib3	Class
11	No	Small	55K	?
12	Yes	Medium	80K	?
13	Yes	Large	110K	?
14	No	Small	95K	?
15	No	Large	67K	?

Test Set



# Examples of Classification Task

- Predicting tumor cells as benign or malignant
- Classifying credit card transactions as legitimate or fraudulent
- Classifying secondary structures of protein as alpha-helix, beta-sheet, or random coil
- Categorizing news stories as finance, weather, entertainment, sports, etc



# Classification Techniques



- Decision Tree based Methods
- Rule-based Methods
- Memory based reasoning
- Neural Networks
- Naïve Bayes and Bayesian Belief Networks
- Support Vector Machines

# Instance Based Classifiers



- Examples:

- Rote-learner

- Memorizes entire training data and performs classification only if attributes of record match one of the training examples exactly

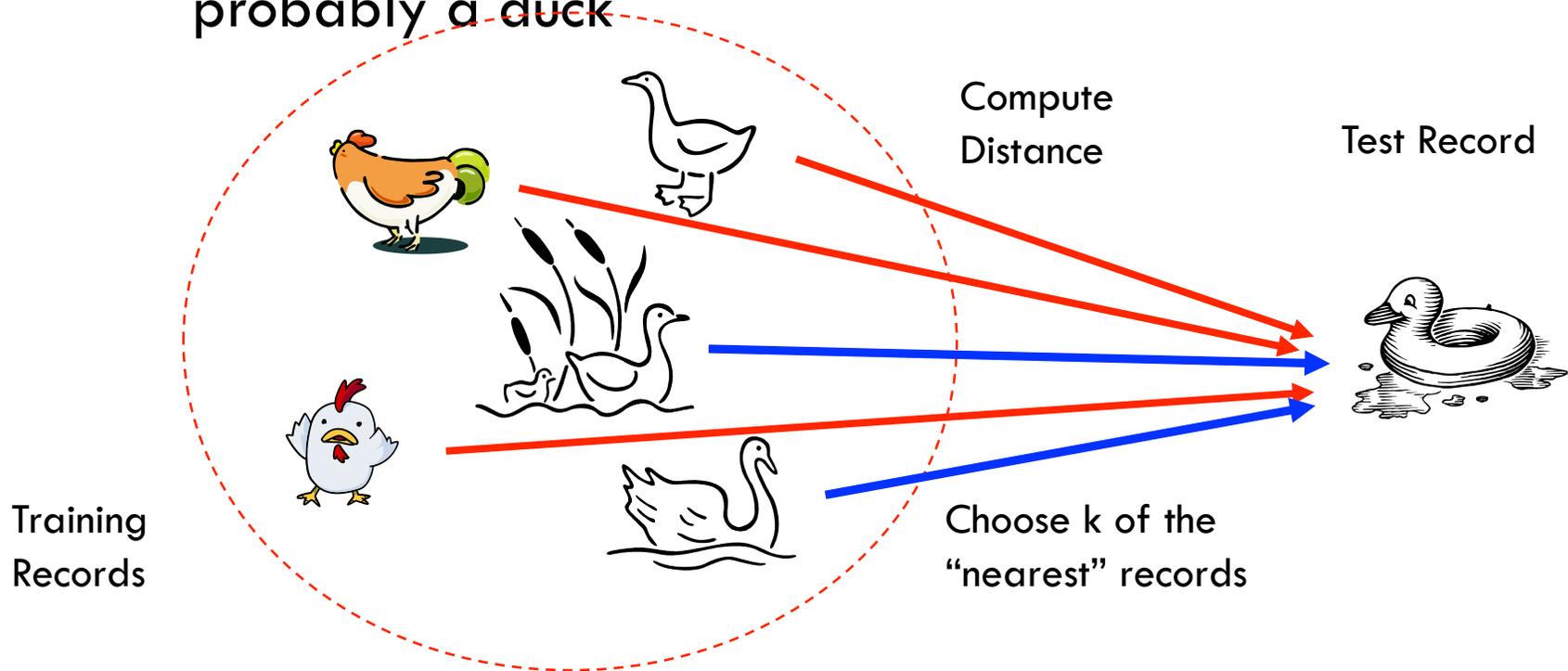
- Nearest neighbor

- Uses  $k$  “closest” points (nearest neighbors) for performing classification

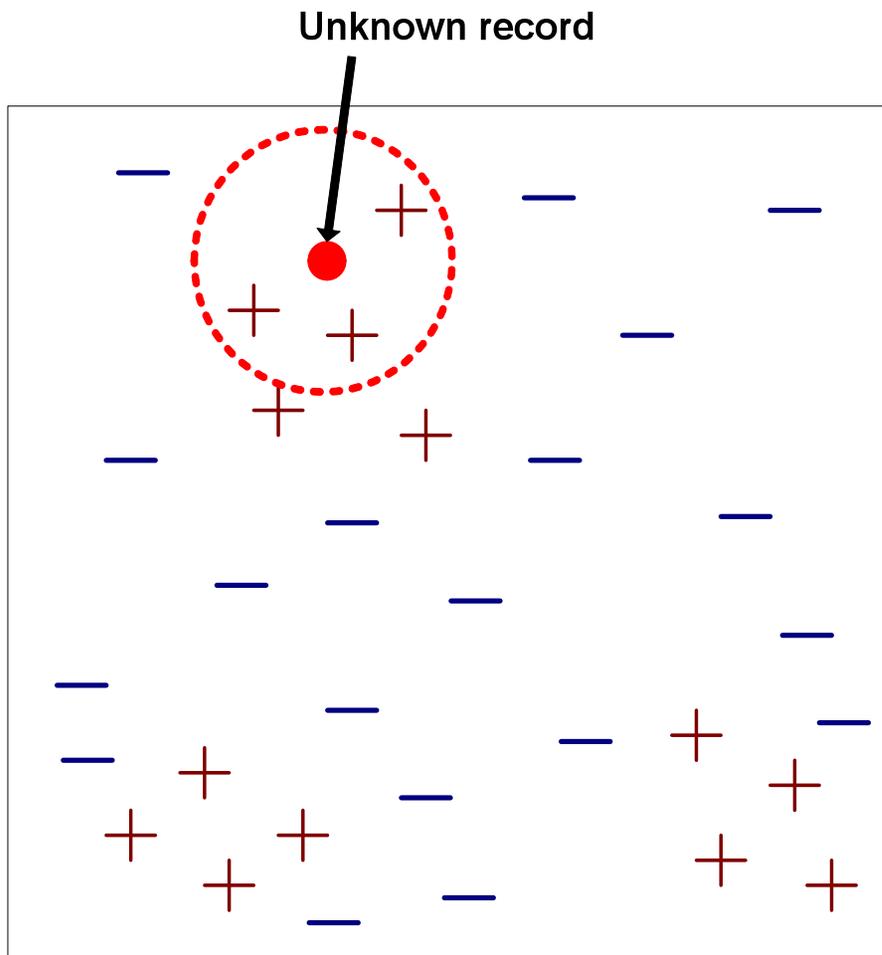
# Nearest Neighbor Classifiers

- Basic idea:

- If it walks like a duck, quacks like a duck, then it's probably a duck

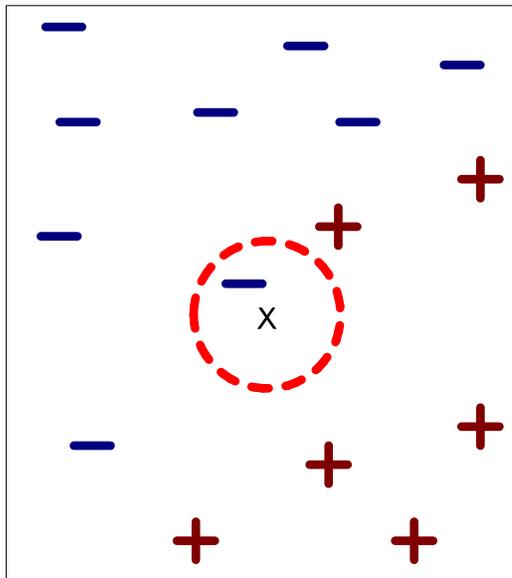


# Nearest-Neighbor Classifiers

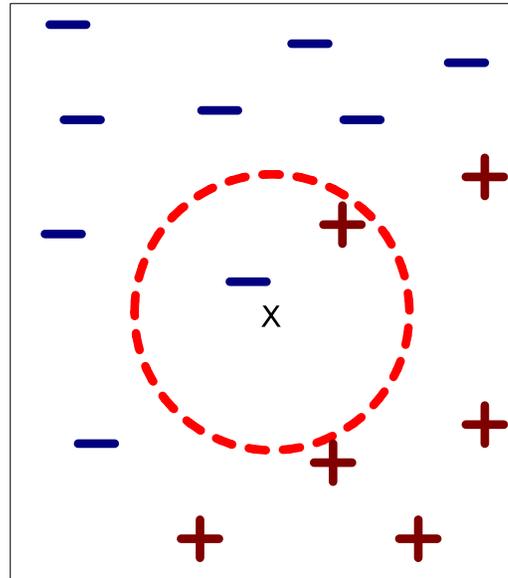


- Requires three things
  - The set of stored records
  - Distance Metric to compute distance between records
  - The value of  $k$ , the number of nearest neighbors to retrieve
- To classify an unknown record:
  - Compute distance to other training records
  - Identify  $k$  nearest neighbors
  - Use class labels of nearest neighbors to determine the class label of unknown record (e.g., by taking majority vote)

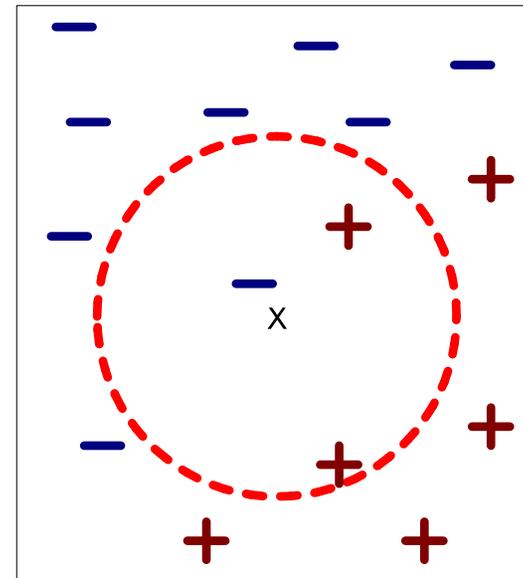
# Definition of Nearest Neighbor



(a) 1-nearest neighbor



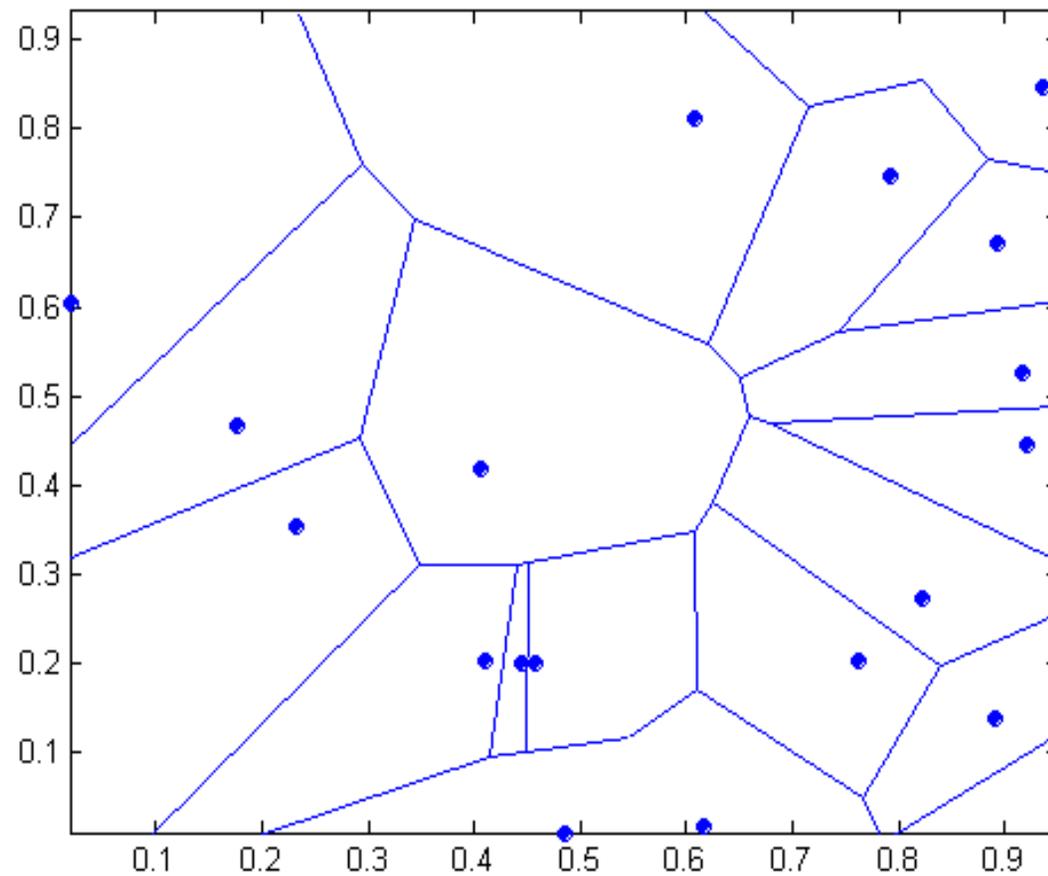
(b) 2-nearest neighbor



(c) 3-nearest neighbor

K-nearest neighbors of a record  $x$  are data points that have the  $k$  smallest distance to  $x$

# 1 nearest-neighbor



Voronoi Diagram

# Nearest Neighbor Classification

- Compute distance between two points:

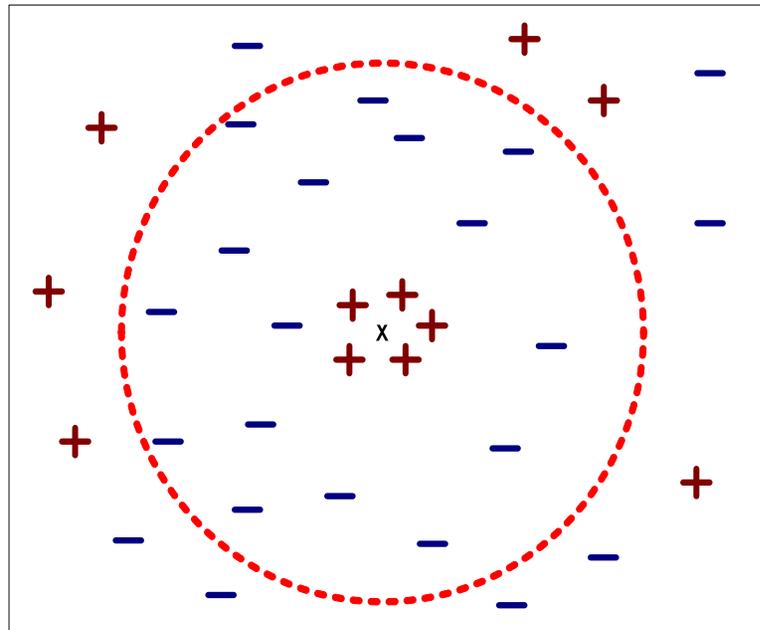
- Euclidean distance

$$d(p, q) = \sqrt{\sum_i (p_i - q_i)^2}$$

- Determine the class from nearest neighbor list
  - take the majority vote of class labels among the k-nearest neighbors
  - Weigh the vote according to distance
    - weight factor,  $w = 1/d^2$

# Nearest Neighbor Classification...

- Choosing the value of  $k$ :
  - If  $k$  is too small, sensitive to noise points
  - If  $k$  is too large, neighborhood may include points from other classes



# Nearest Neighbor Classification...



- Scaling issues

- Attributes may have to be scaled to prevent distance measures from being dominated by one of the attributes

- Example:

- height of a person may vary from 1.5m to 1.8m
    - weight of a person may vary from 90lb to 300lb
    - income of a person may vary from \$10K to \$1M

# Evaluating Performance of Classifier



- Instances are partitioned into TRAIN and TEST sets
- TRAIN set is used to build the model
- TEST is used to evaluate the classifier
- Methods for creating TRAIN and TEST sets
  - Holdout
  - Random Sub-sampling
  - Cross-Validation
  - Bootstrap

# Classifier Accuracy



- Accuracy
  - Number of instances correctly classified
- Problematic for unbalanced classes
  - 990 Class A
  - 10 Class B
  - Classifier always predicts A has accuracy = 99%

# Unbalanced Classes

		actual value		total
		$p$	$n$	
prediction outcome	$p'$	True Positive	False Positive	$P'$
	$n'$	False Negative	True Negative	$N'$
total		$P$	$N$	

- Precision =  $TP/P'$

- Recall =  $TP/P$

- F1

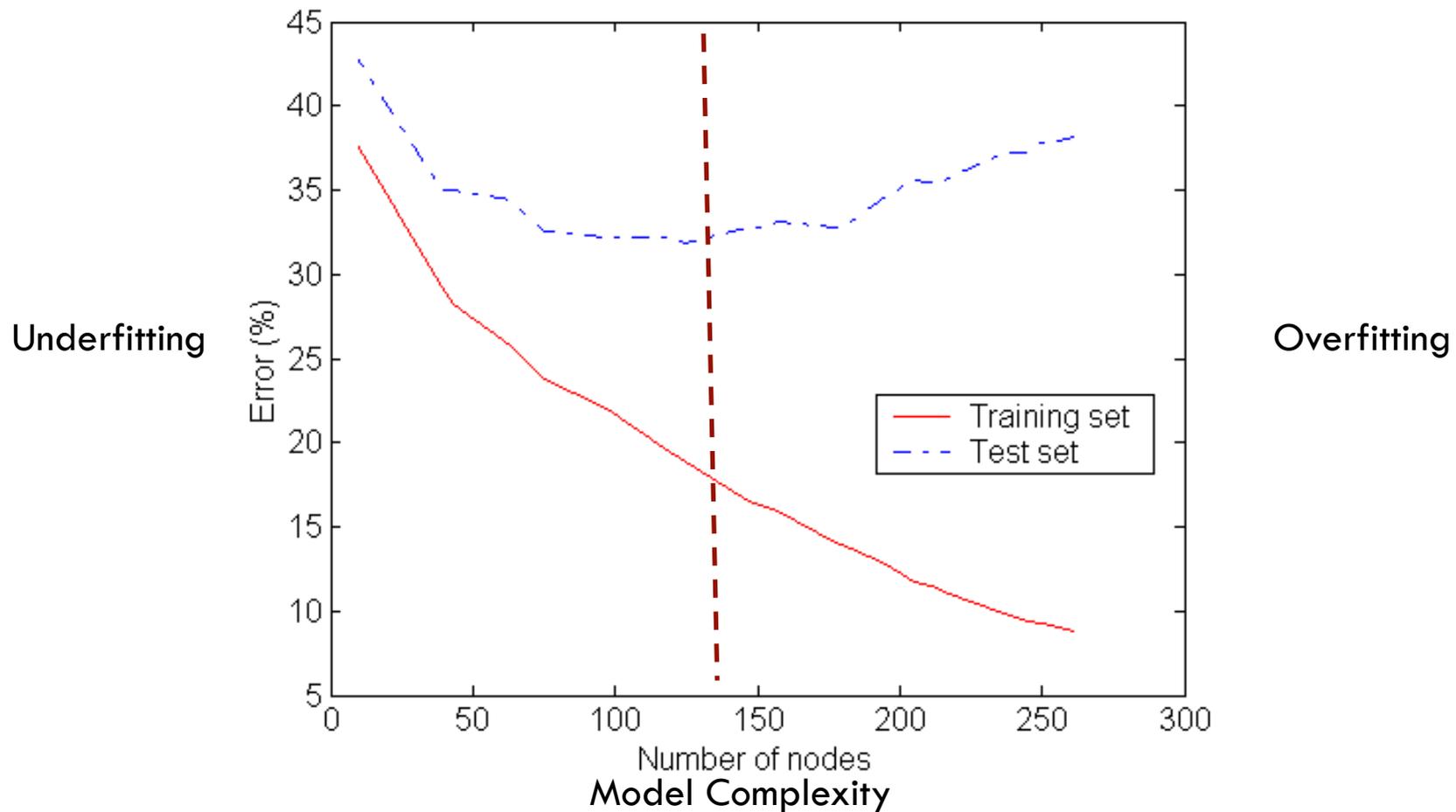
- Sensitivity =  $TP/P$

- Specificity =  $TN/N$

- ROC

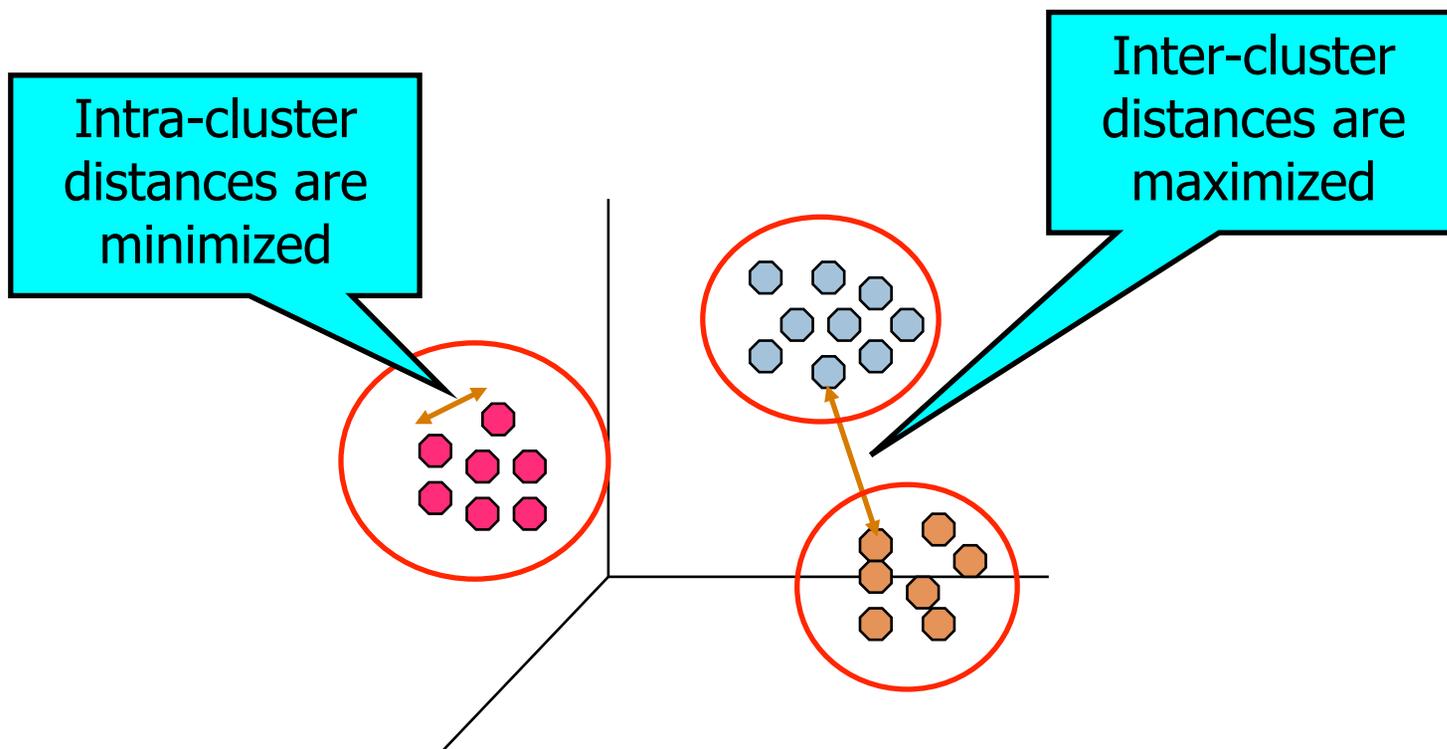
(receiver operating characteristic)

# Model Generalization

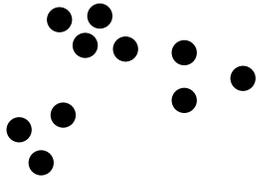


# What is Cluster Analysis?

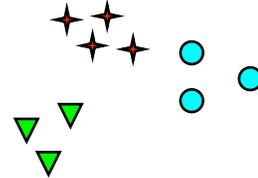
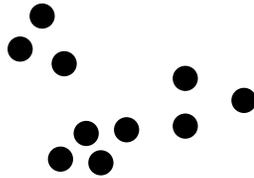
- Finding groups of objects such that the objects in a group will be similar (or related) to one another and different from (or unrelated to) the objects in other groups



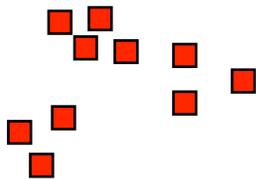
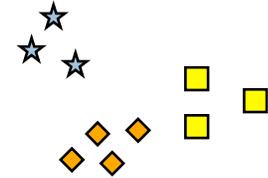
# Notion of a Cluster can be Ambiguous



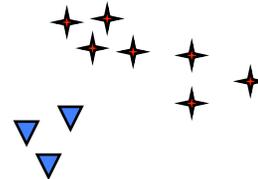
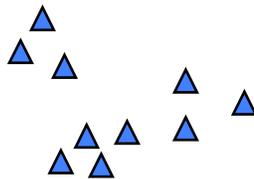
How many clusters?



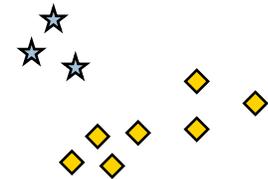
Six Clusters



Two Clusters



Four Clusters



# Types of Clusterings



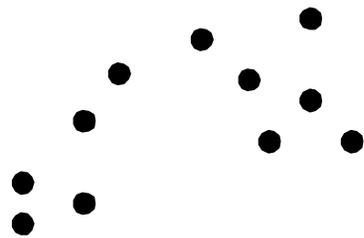
- Partitional Clustering

- A division data objects into non-overlapping subsets (clusters) such that each data object is in exactly one subset

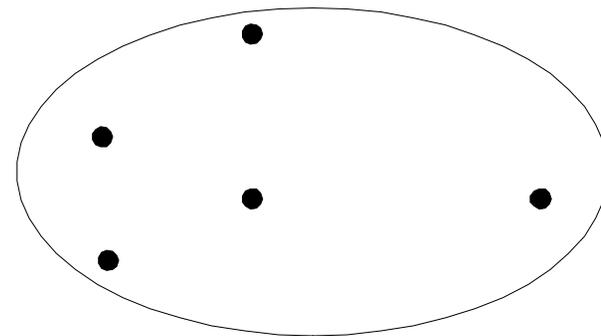
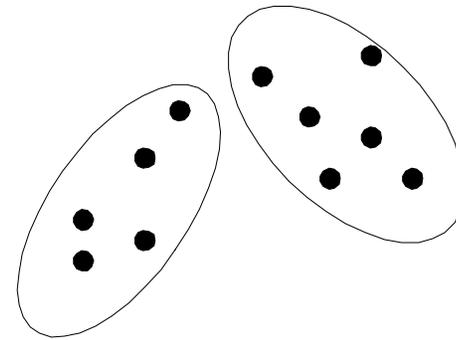
- Hierarchical clustering

- A set of nested clusters organized as a hierarchical tree

# Partitional Clustering

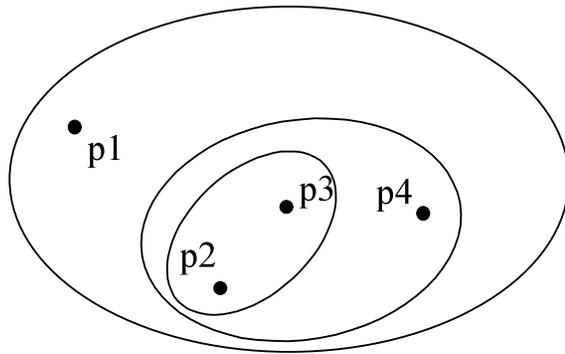


Original Points

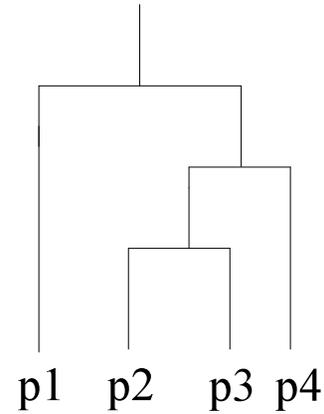


A Partitional Clustering

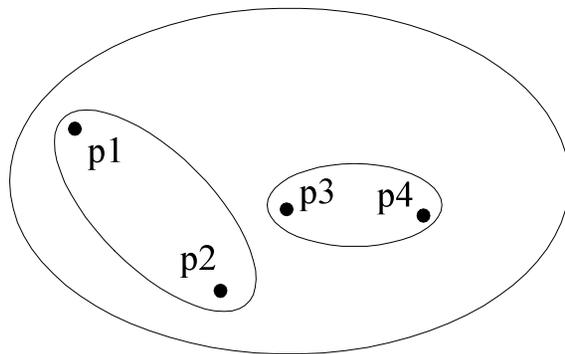
# Hierarchical Clustering



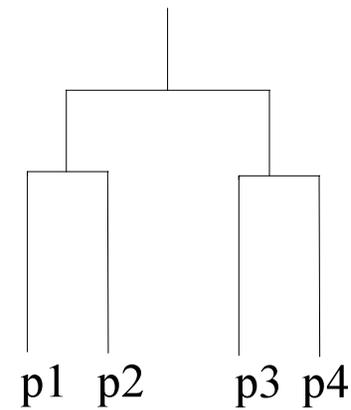
Traditional Hierarchical Clustering



Traditional Dendrogram



Non-traditional Hierarchical Clustering



Non-traditional Dendrogram

# K-means Clustering

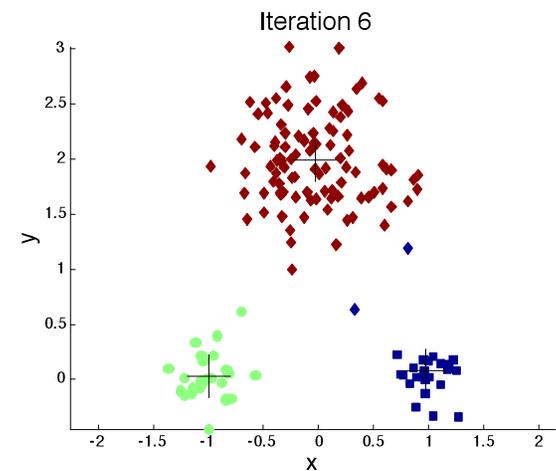
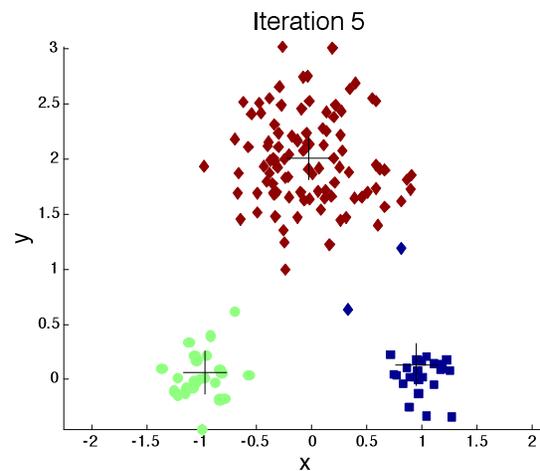
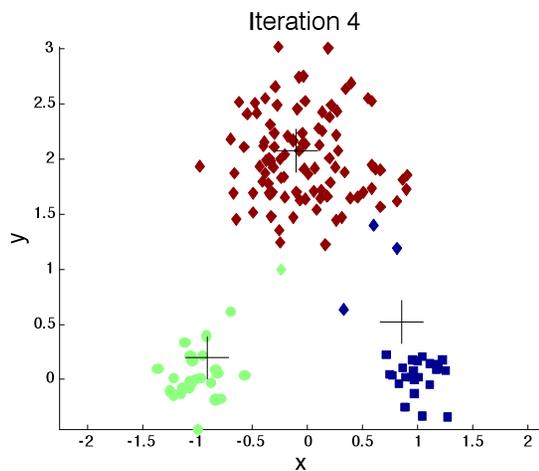
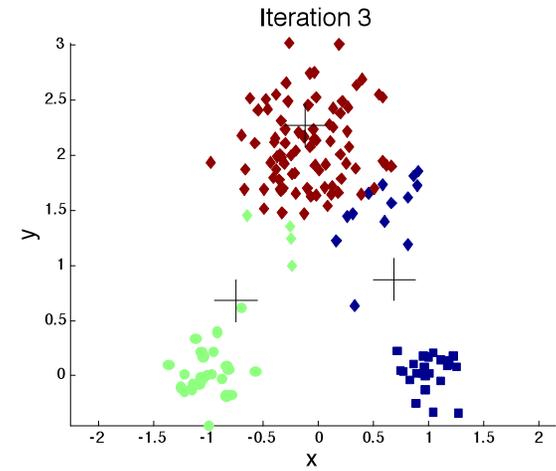
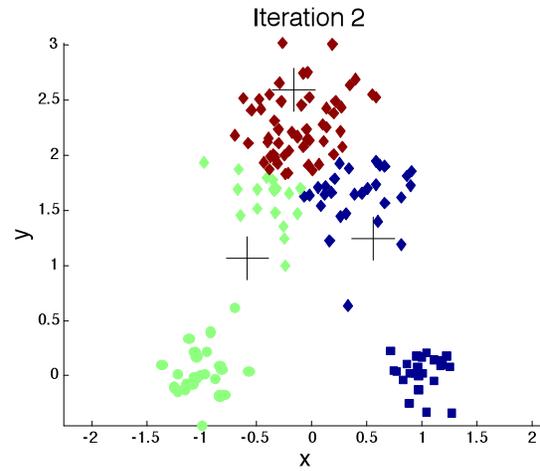
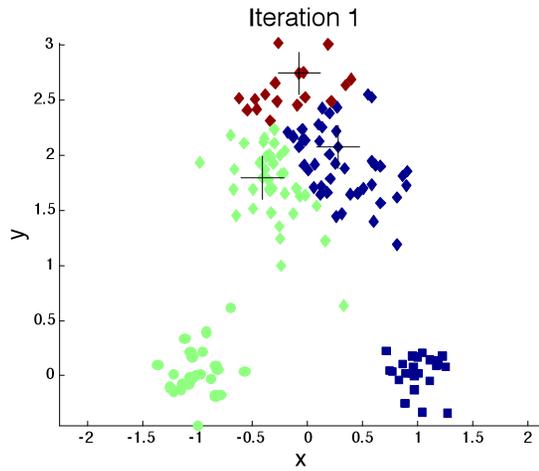
- Partitional clustering approach
- Prototype-based
- Each cluster is associated with a **centroid** (center point)
- Each point is assigned to the cluster with the closest centroid
- Number of clusters,  $K$ , must be specified

- 
- 1: Select  $K$  points as the initial centroids.
  - 2: **repeat**
  - 3: Form  $K$  clusters by assigning all points to the closest centroid.
  - 4: Recompute the centroid of each cluster.
  - 5: **until** The centroids don't change
-

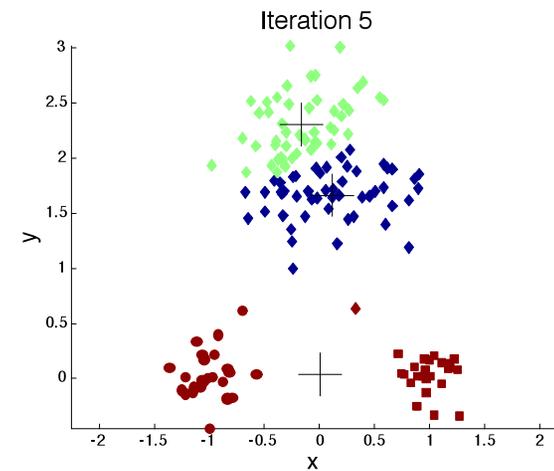
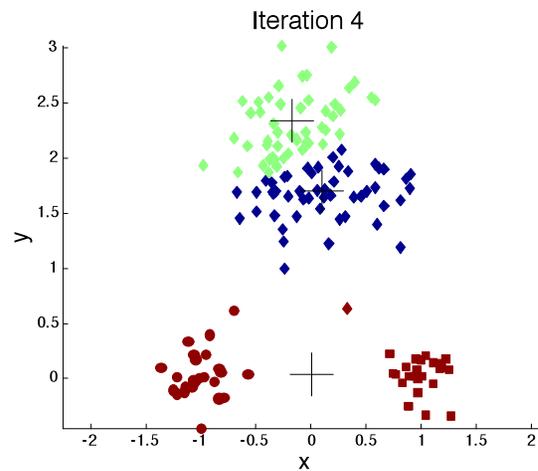
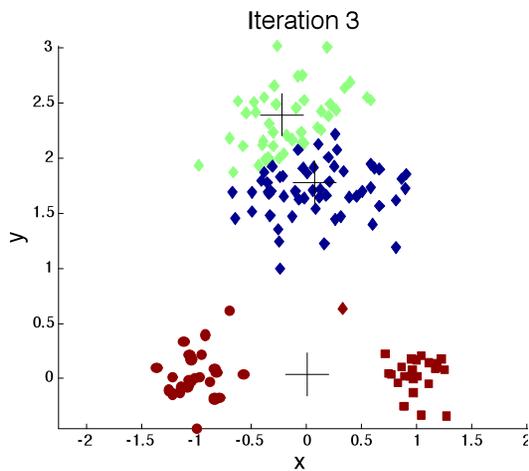
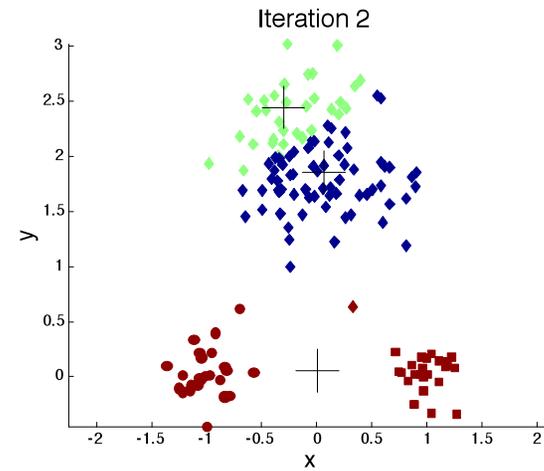
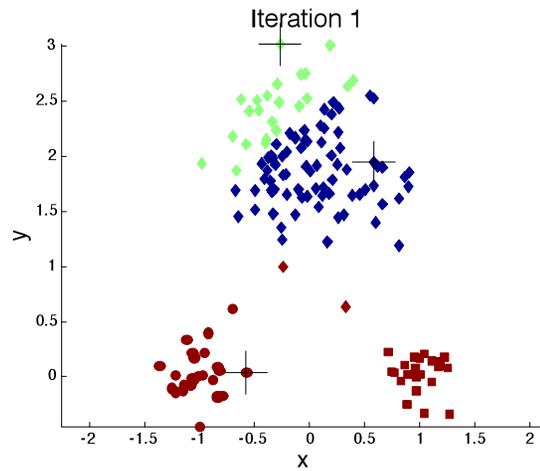
# K-means Clustering – Details

- Initial centroids are often chosen randomly.
  - ▣ Clusters produced vary from one run to another.
- The centroid is (typically) the mean of the points in the cluster.
- ‘Closeness’ is measured by Euclidean distance, cosine similarity, correlation, etc.
- K-means will converge for common similarity measures mentioned above.
- Most of the convergence happens in the first few iterations.
  - ▣ Often the stopping condition is changed to ‘Until relatively few points change clusters’
- Complexity is  $O( n * K * I * d )$ 
  - ▣  $n$  = number of points,  $K$  = number of clusters,  
 $I$  = number of iterations,  $d$  = number of attributes

# Importance of Choosing Initial Centroids



# Importance of Choosing Initial Centroids ...



# Limitations of K-means



- K-means has problems when clusters are of differing
  - Sizes
  - Densities
  - Non-globular shapes
- K-means has problems when the data contains outliers.

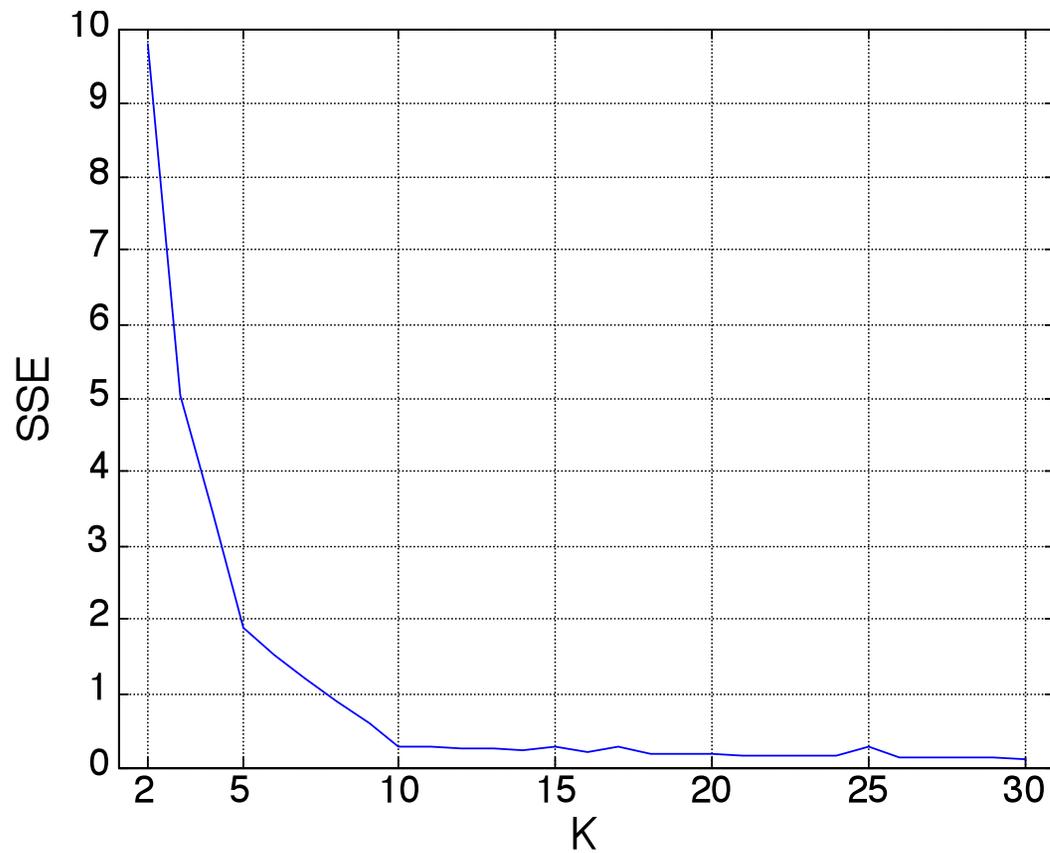
# Evaluating K-means Clusters

- Most common measure is Sum of Squared Error (SSE)
  - ▣ For each point, the error is the distance to the nearest cluster
  - ▣ To get SSE, we square these errors and sum them.

$$SSE = \sum_{i=1}^K \sum_{x \in C_i} dist(c_i, x)^2$$

- ▣  $x$  is a data point in cluster  $C_i$  and  $c_i$  is the representative point for cluster  $C_i$
- ▣ Given two sets of clusters, we can choose the one with the smallest error
- ▣ One easy way to reduce SSE is to increase  $K$ , the number of clusters

# Relative Index



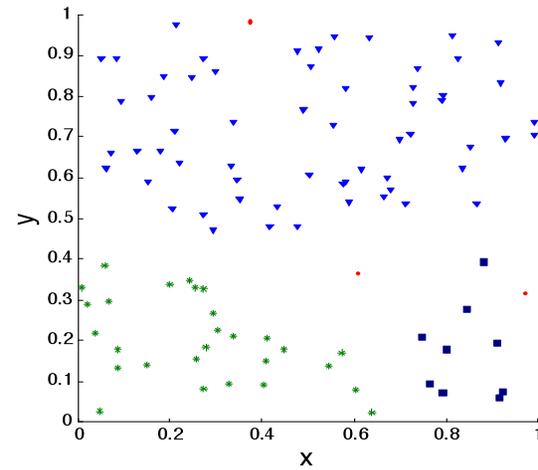
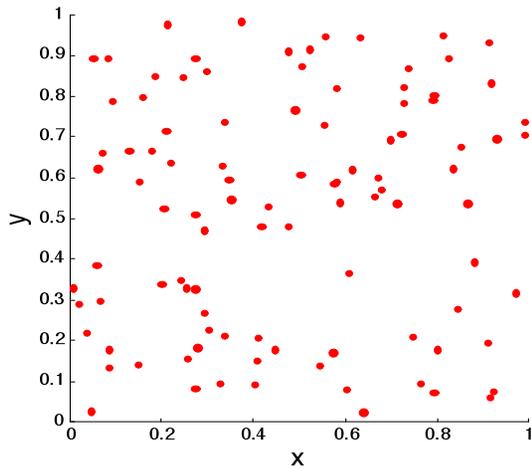
K-Means: SSE vs Number of Clusters

# Cluster Validity

- For supervised classification we have a variety of measures to evaluate how good our model is
  - ▣ Accuracy, precision, recall
  
- For cluster analysis, the analogous question is how to evaluate the “goodness” of the resulting clusters?
  
- But “clusters are in the eye of the beholder”!
  
- Then why do we want to evaluate them?
  - ▣ To avoid finding patterns in noise
  - ▣ To compare clustering algorithms
  - ▣ To compare two sets of clusters
  - ▣ To compare two clusters

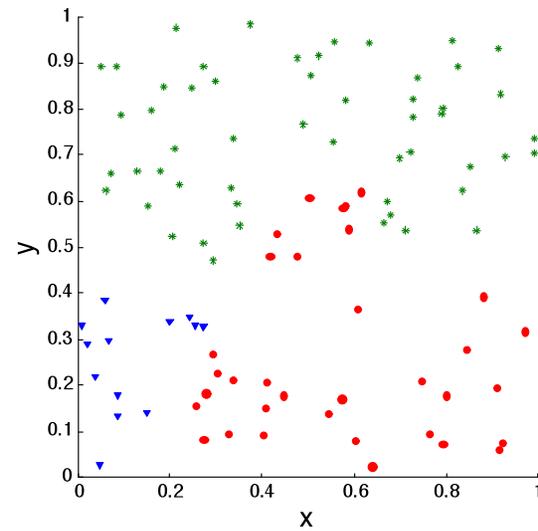
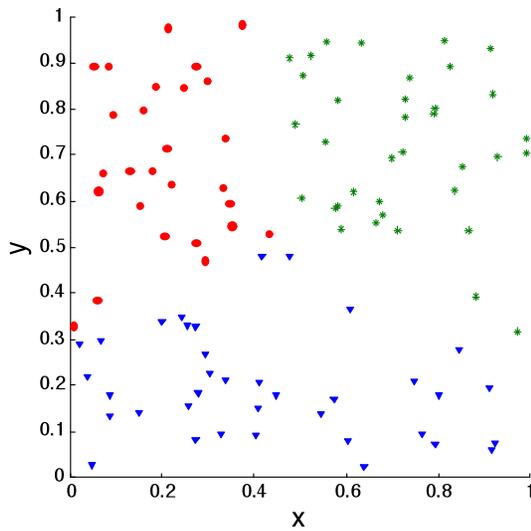
# Clusters found in Random Data

Random  
Points



DBSCAN

K-means



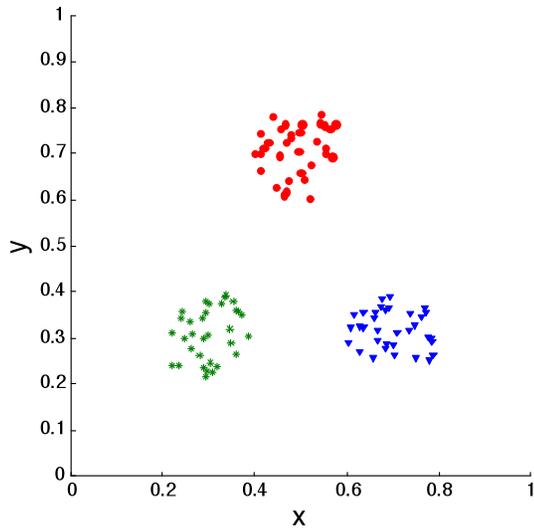
Complete  
Link

# Measures of Cluster Validity

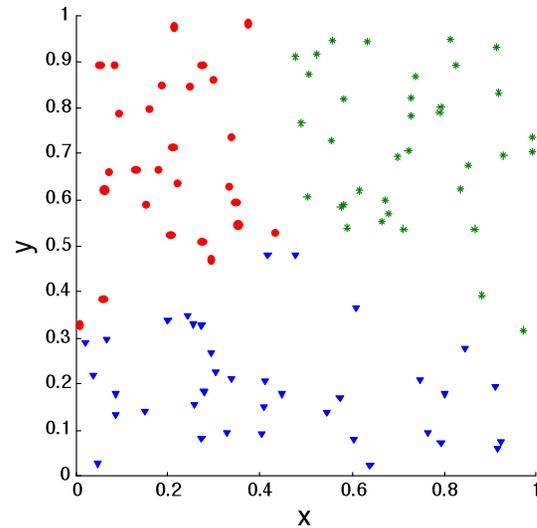
- Numerical measures that are applied to judge various aspects of cluster validity, are classified into the following three types.
  - **External Index:** Used to measure the extent to which cluster labels match externally supplied class labels.
    - Entropy
  - **Internal Index:** Used to measure the goodness of a clustering structure *without* respect to external information.
    - Sum of Squared Error (SSE)
  - **Relative Index:** Used to compare two different clusterings or clusters.
    - Often an external or internal index is used for this function, e.g., SSE or entropy

# Measuring Cluster Validity Via Correlation

- Correlation of incidence and proximity matrices for the K-means clusterings of the following two data sets.



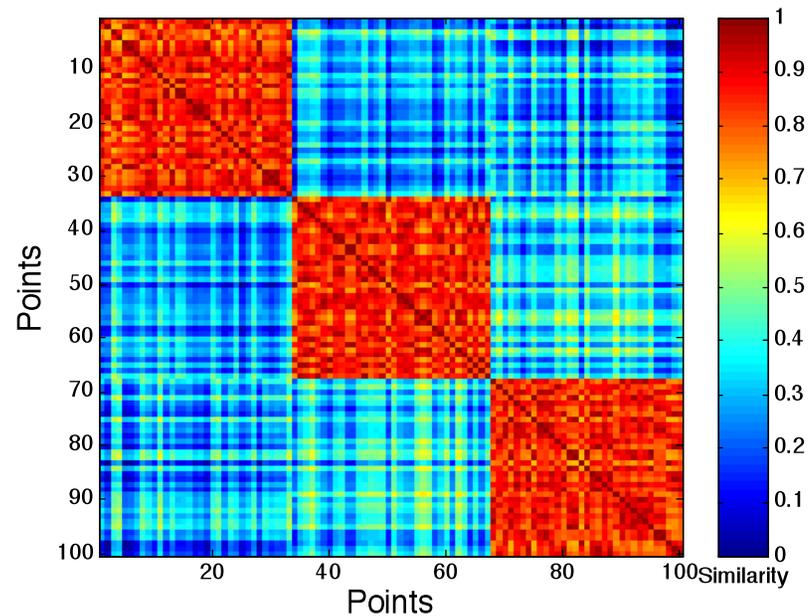
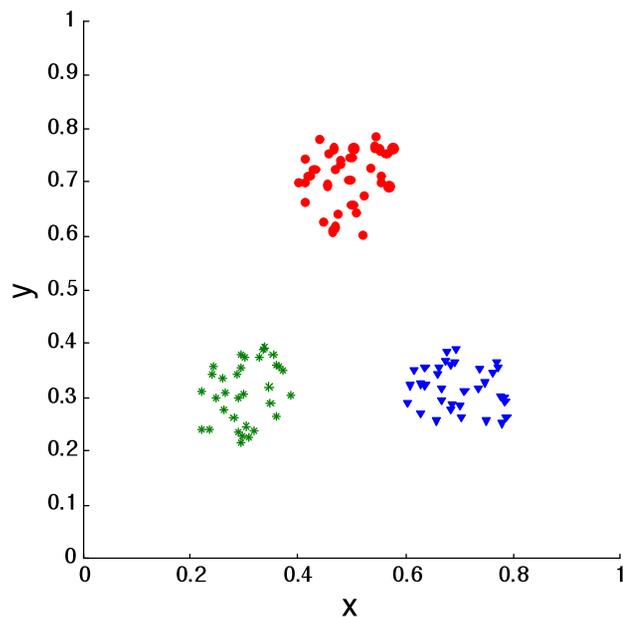
Corr = -0.9235



Corr = -0.5810

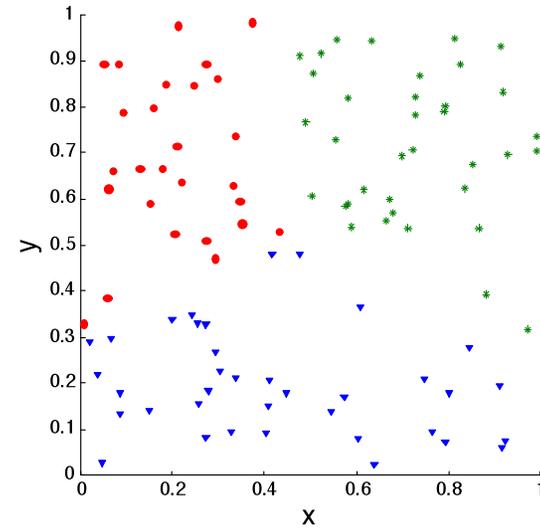
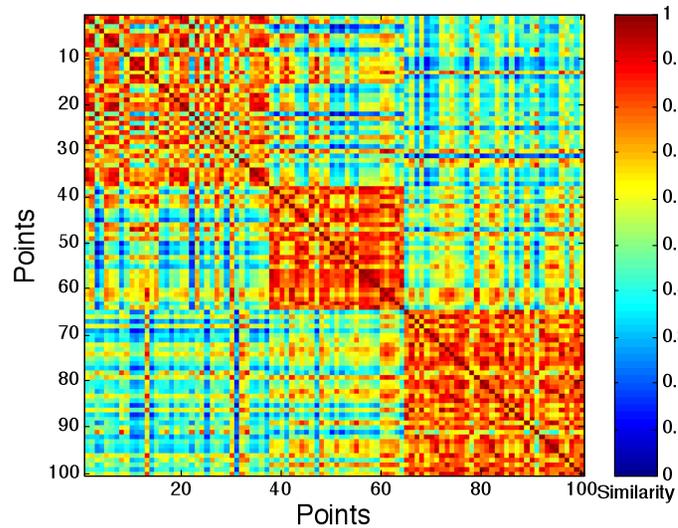
# Using Similarity Matrix for Cluster Validation

- Order the similarity matrix with respect to cluster labels and inspect visually.



# Using Similarity Matrix for Cluster Validation

- Clusters in random data are not so crisp



K-means