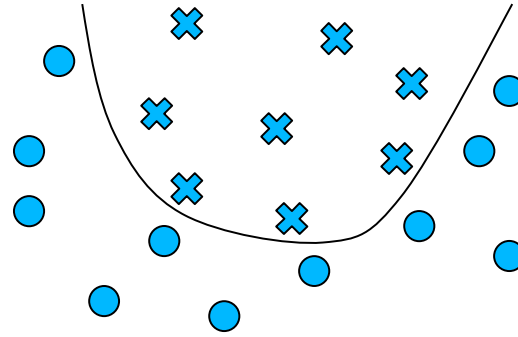# Pattern Recognition (Örüntü Tanıma)

## Basic Approaches and Techniques

### Neşe Yalabık

Hasat Projesi Eğitim Programı

ODTÜ, 31 Ağustos 2010

# *Outline*

Pattern Recognition: Definitions and Objectives

Bayes Classifiers

Estimation of Parameters

Features and Feature Selection

K –Nearest Neighbor Classifiers

Linear Discriminant Function Classifiers

Tree Classifiers

Clustering

Training and Performance Testing in PR

References

# *Pattern Recognition: Definitions*

**Pattern Recognition(PR):** The process of machine perception for an automatic labeling of an object or an event into one of the predefined categories.

**Pattern Classification:** Final step in a PR system

_____

We human beings do pattern recognition everyday.
We "**recognize**" and **classify** many things,
even if it is corrupted by **noise**, **distorted** and **variable**.

Classification is the result of **recognition: learning, categorization, generalization**

A problem is a PR problem only if it involves '**statistical variation**'

# *Example*

We see here that all 9's are different from each other and 9's and 4's can easily be mixed

| | | | | |
|---|---|---|---|---|
| 1 | 9 | 9 | 3 | Recognized as 1393 |
| 1 | 9 | 9 | 7 | Recognized as 1937 |
| 1 | 9 | 9 | 4 | Recognized as 1434 |
| 1 | 0 | 6 | 8 | Recognized as 1060 |
| 1 | 9 | 9 | 4 | Recognized as 1394 |
| 1 | 9 | 4 | 5 | Recognized as 1995 |
| 1 | 9 | 4 | 8 | Recognized as 1940 |
| 1 | 9 | 9 | 0 | Recognized as 1930 |
| 1 | 9 | 4 | 5 | Recognized as 1995 |
| 1 | 9 | 7 | 3 | Recognized as 1573 |
| 1 | 9 | 8 | 3 | Recognized as 1583 |
| 1 | 9 | 9 | 1 | Recognized as 1951 |

# Example Applications of Pattern Recognition

- Reading hand-written text to classify it into letters and words

- Analyzing fingerprints to find the owner

- Recognizing the faces of people to name them

- Finding buildings in a satellite image

- Naming a gun from its bullet mark (Ballistics)

- Identifying different objects on a conveyor belt

- Analyzing test results in decision support for any illness

# *Pattern Recognition: Definitions*

A Pattern Recognition System consists of the following parts:

    **Pre-processing and Feature Extraction**

    **Learning**

    **Classification**

    **Post- processing**

# *Pattern Recognition: Definitions*

**"Pre- processing and Feature Extraction"** Converts 'data' to 'features'

**"Data"** raw data taken as samples

**"Feature"** a discriminating, easily measurable characteristics of our data.

**"Feature Vector":** A set of variables that represent different features

**"Feature Space" :** is defined by a feature vector

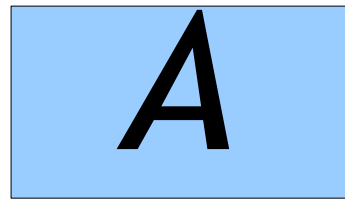# Pattern Recognition: Definitions

## Pre-processing and Pattern Extraction

A ll operations over raw data (such as a remotely sensed image) to enhance and process it , leading to extraction of features .

Includes enhancement, edge extraction, segmentation etc.

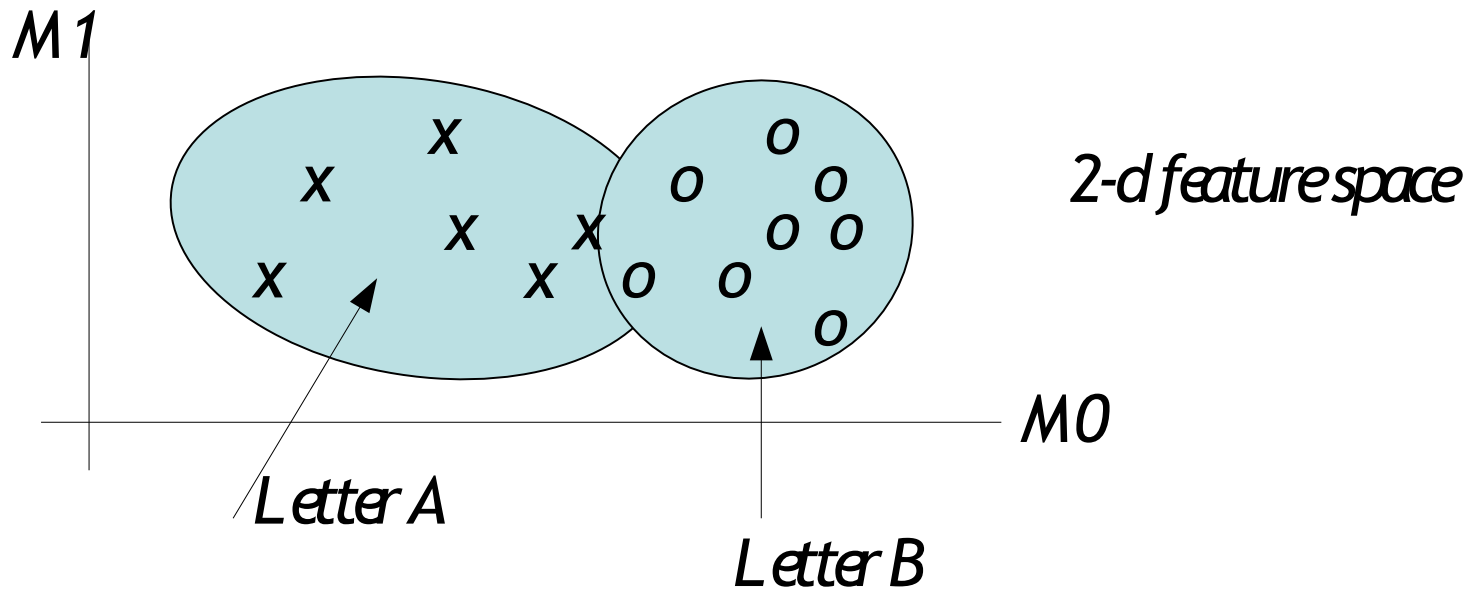Pattern extraction results with a feature vector X

# *Example*
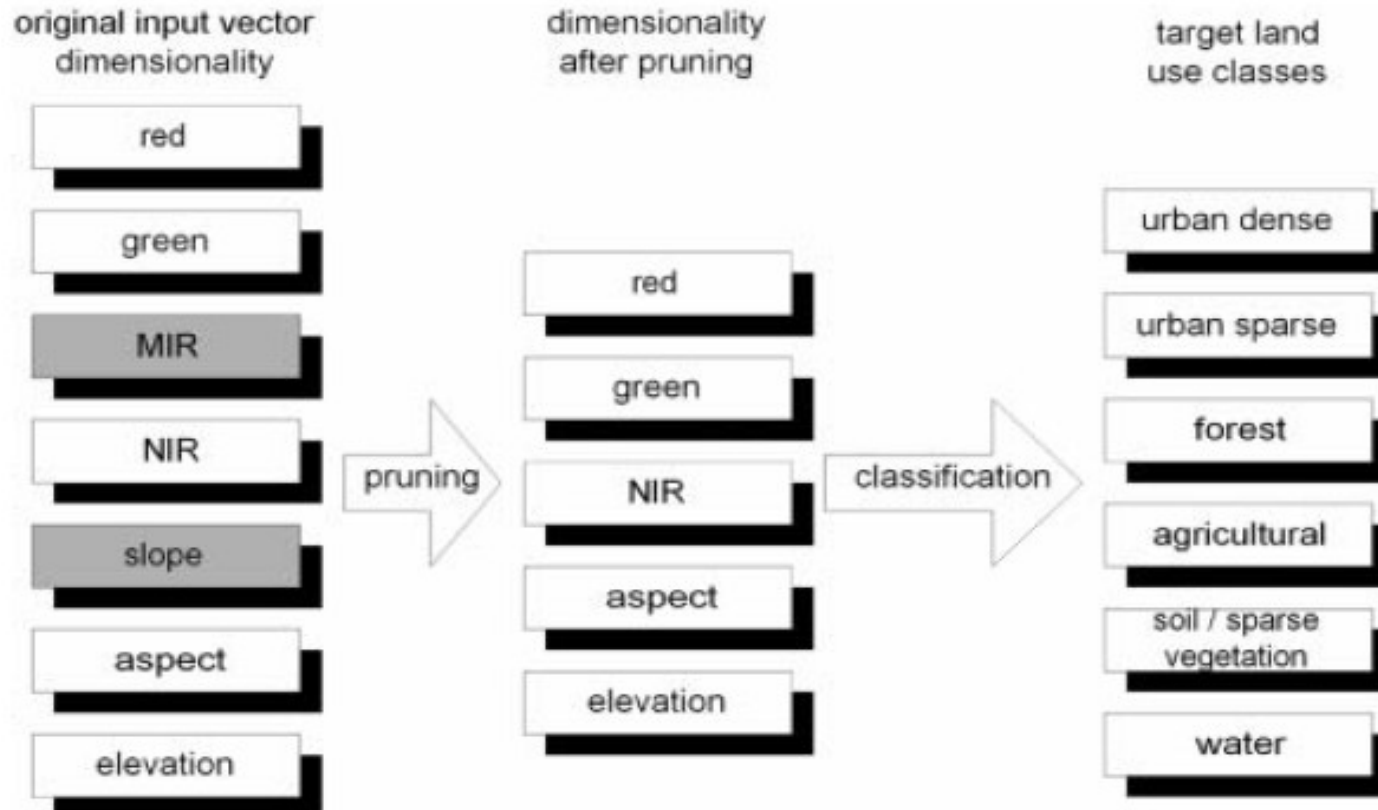
Consider recognition of handwritten characters:

A

$$[M_0, M_1, ..., M_k]$$

Raw data: Bitmap       Features: Moment invariants

*M1*

x   x   x   x   x   x

o   o   o   o   o   o   o   o

*2-d feature space*

*M0*

*Letter A*

*Letter B*

# *Example*

## Clasification of land use from multispectral data

## Satellite image classification using granular neural networks

D. Stathakis[ab]; A. Vasilakos[a]
[a] Department of Planning and Regional Development, University of Thessaly, Pedion Areos 38334, Greece [b] EC Joint Research Centre, IPSC, MARS-FOOD, 21020 (VA), Italy

# *Pattern Recognition: Definitions*

**Learning:** Devising a classifier from collected samples with or without known labels (categories)

**"Learning samples"** Large data sets to be used in training, or estimating parameters, etc. They may be labeled or not.

Given the learning data set with known labels:**supervised learning;**Unknown labels: **unsupervised learning and clustering**

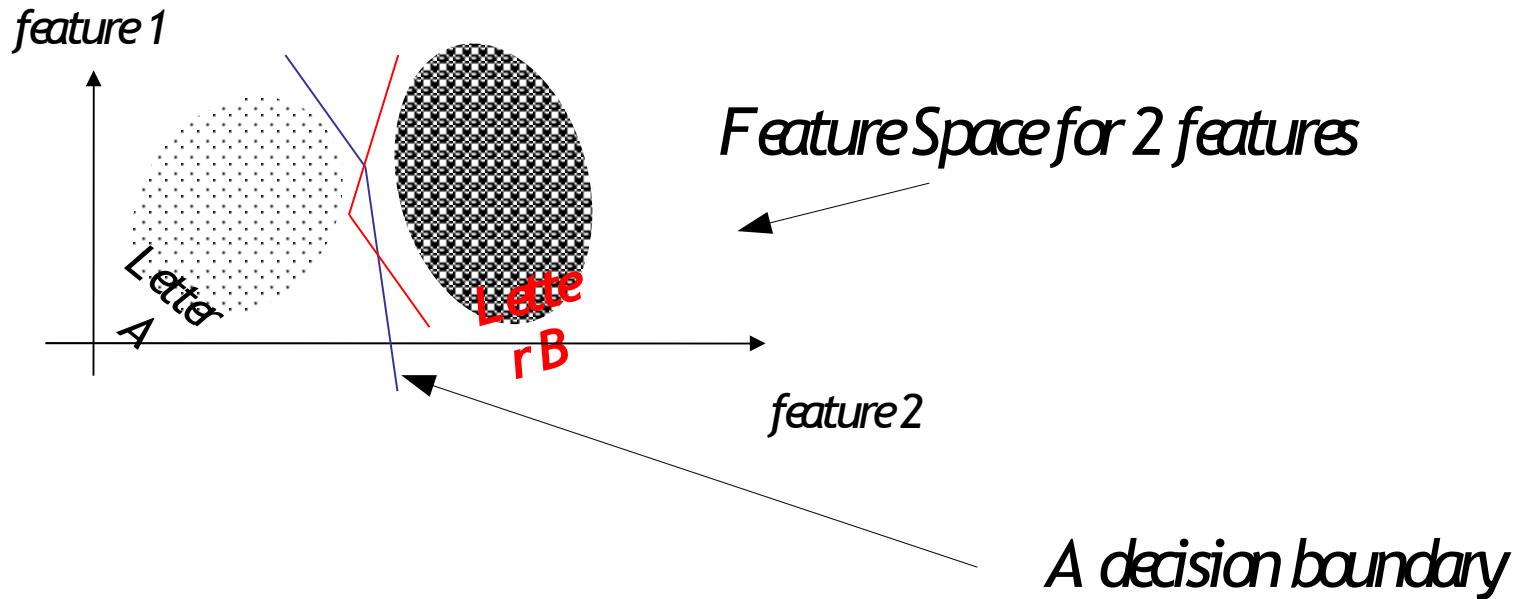**"Test Samples"** used in testing the classifier performance.

**"Result"** a decision on the category sample belongs.

**"Performance"** How well a classifier correctly recognizes test samples
**"Correct Classification ratio"** ratio of correctly classified samples to all test samples

# Pattern Recognition: Definitions

## Classification

feature 1

Letter A

Letter B

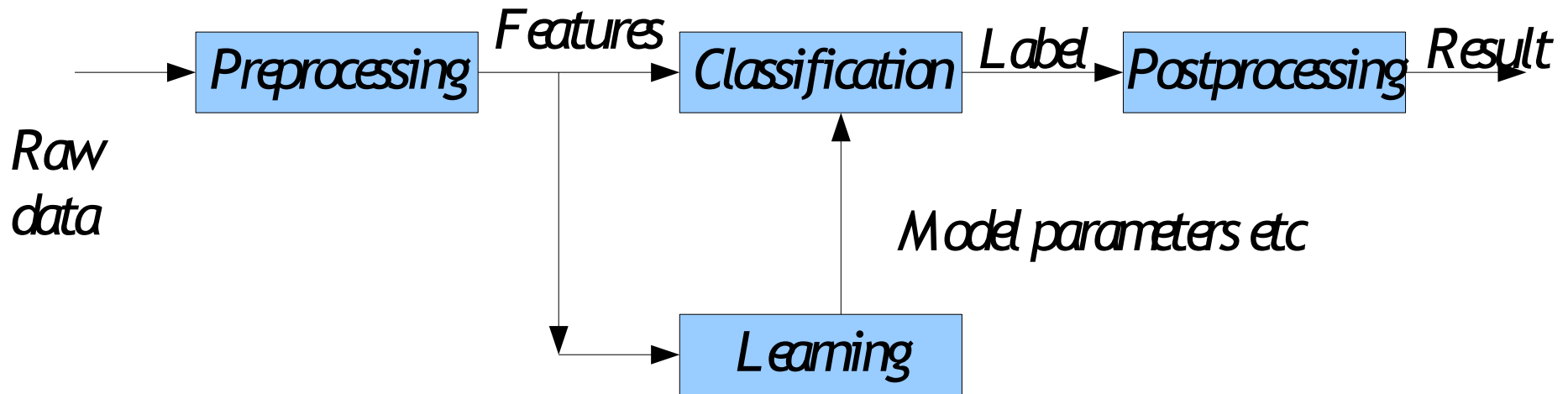Feature Space for 2 features

feature 2

A decision boundary

How do we separate A 's from B 's?
Form a decision boundary
Classify the sample to the side it falls

# Pattern Recognition: Definitions

**Post- processing**: Domain knowledge may be incorparated to correct mistakes, such as using language to correct letter classifiers

## A Pattern Recognition System

Raw data → Preprocessing → *Features* → Classification → *Label* → Postprocessing → *Result*

Preprocessing → Learning

Learning → *Model parameters etc* → Classification

# *Objective in PR*

**Performance criteria: Minimize the average error** (at least as good as a human being )

**Minimize the risk:** wrong decision could be more risky in some cases such as medical diagnosis

**Why automise?** Obvious reason: save from time and effort

(Ex: consensus forms: enter 100 million records into electronic medium ).

**How do machines solve it:** Many different approaches in history

**Statistical Pattern Recognition:** relies on statistics of collected data

**Structural Pattern Recognition:** tries to discover the structure inherent in data

(ex: may assume letters are composed of strokes etc )

# Statistical Approach to P.R

$$X = [X_1, X_2, ..., X_d]$$

*Dimension of the feature space:* $\quad d$

*Set of different states of nature:* $\quad \{\omega_1, \omega_2, ..., \omega_c\}$

*Categories:* $c$

*find*

$$R_i \quad R_i \cap R_j = \emptyset \quad uR_i = R^d$$

*set of possible actions (decisions):* $\quad \{\alpha_1, \alpha_2, ..., \alpha_a\}$
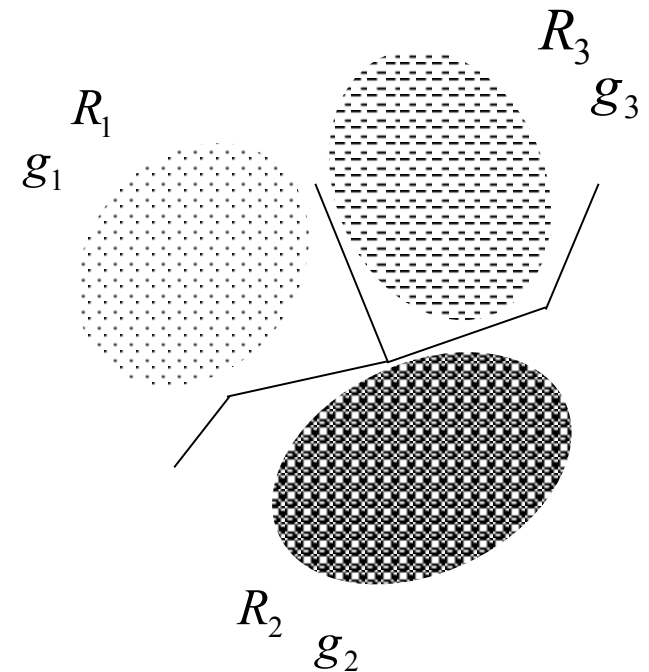
*Here, a decision might include a 'reject option'*

<u>*A Discriminant Function*</u> $\qquad\qquad\qquad g_i(X) \qquad\qquad g_i(X) \geq g_j(X)$

*in region* $R_i$ *; decision rule :* $\quad$ *if* $\alpha_k \qquad g_k(X) > g_j(X) \qquad\qquad 1 \leq i \leq c$

$R_3$

$g_3$

$R_1$

$g_1$

$R_2$

$g_2$

# A Pattern Classifier



So our aim now will be to define these functions $g_1, g_2, ..., g_c$ to *minimize* or *optimize* a criterion.

# *Bayes Classifiers*

A Parametric approach which assumes that the feature vectors are random variables with known probability distributions.

'Bayes Decision Theory' is used or minimum –error– minimum risk pattern classifier design.

It is assumed that if a sample X is drawn from category w i, it is a random variable represented with a multivariate probability density function called

'Class–conditional density function'

$$P(X \mid Wi)$$

We also know a-priori probability $P(\omega_i)$

Then, we can talk about a decision rule
that minimizes the probability of error.

Suppose we have the observation  X
This observation is going to change a-priori assumption
to a-posteriori probability:
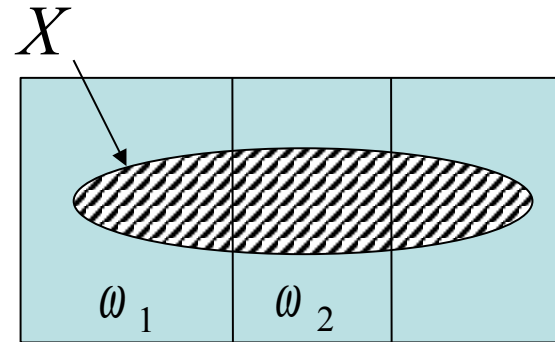
$$P(\omega_i|X)$$

which can be found by the Bayes Rule

$$P(\omega_i \mid X) = P(\omega_i, X) / P(X)$$

$$= \frac{P(X \mid \omega_i) . P(\omega_i)}{P(X)}$$
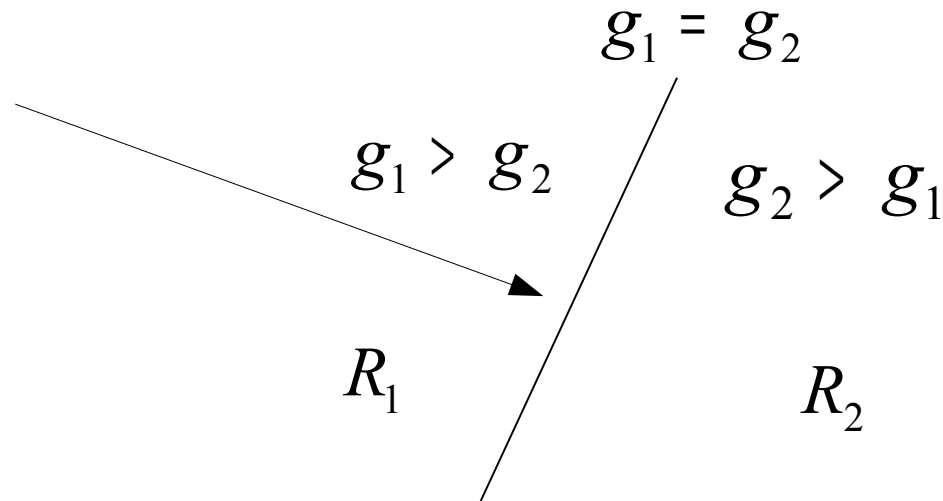
$P(X)$ *can be found by Total Probability Rule:*

*When*

$$P(X) = \sum_{i=1}^{c} P(\omega_i, X)$$



**Decision Rule:** Choose the category with highest a-posteriori probability, calculated as above, using Bayes Rule.

*then,* $g_i(X) = P(\omega_i | X)$

*Decision boundary:*

$g_1 = g_2$

$g_1 > g_2$     $g_2 > g_1$

$R_1$     $R_2$

*or in general, decision boundaries are where:*

$$g_i(X) = g_j(X)$$

*between regions* $R_i$ *and* $R_j$

*Single feature –  decision boundary –  point*
*2 features –   curve (quadratic for gaussian distribution)*
*3 features –      surface*
*More than 3 –       hypersurface*

$$g_i(X) = P(X|\omega_i).P(\omega_i)$$

$$gi(X) = \frac{P(X|\omega_i).P(\omega_i)}{P(X)}$$

*Sometimes, it is easier to work with logarithms*

$$g_i(X) = \log[P(X|\omega_i).P(\omega_i)]$$

$$g_i(X) = \log P(X|\omega_i) + \log P(\omega_i)$$

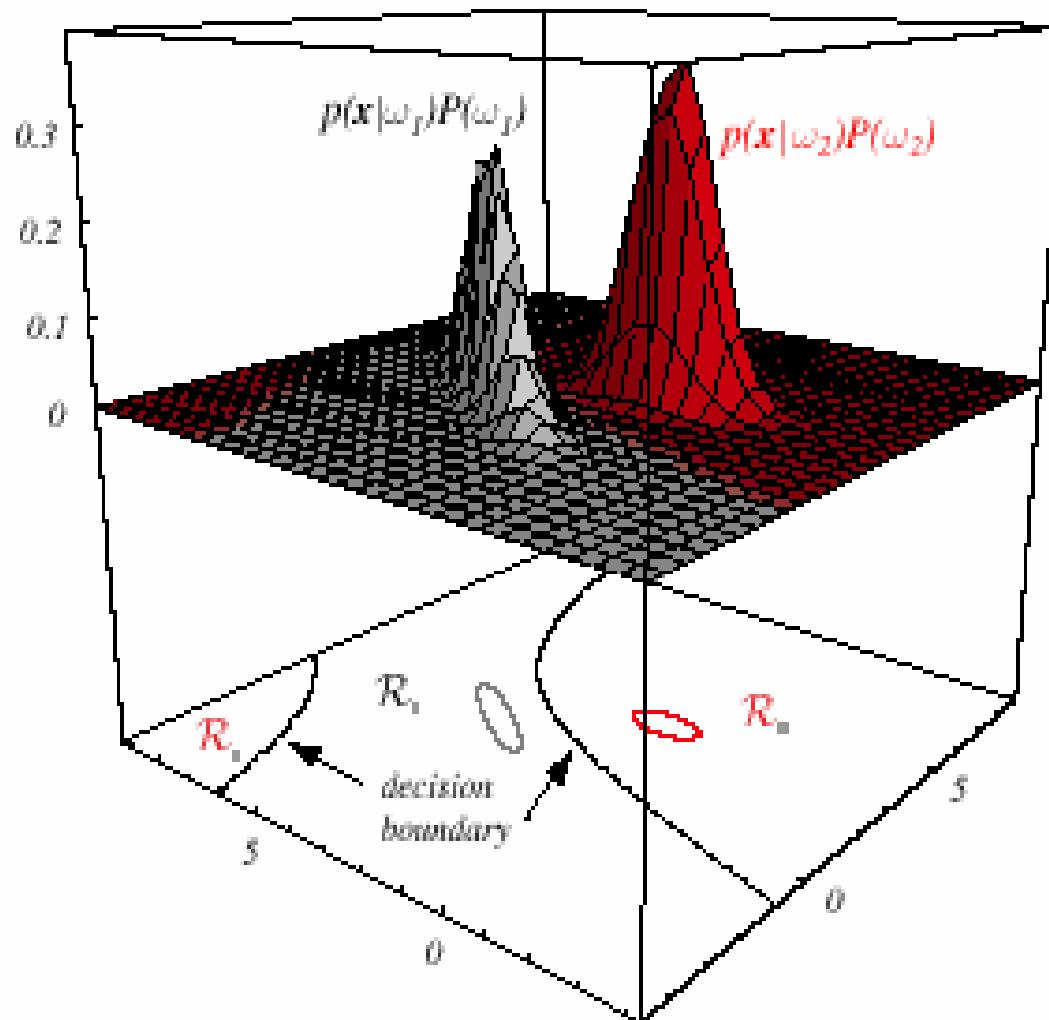*Since logarithmic function is a monotonically increasing function, log fn.will give the same result.*

**FIGURE 2.6.** In this two-dimensional two-category classifier, the probability densities are Gaussian, the decision boundary consists of two hyperbolas, and thus the decision region $\mathcal{R}_2$ is not simply connected. The ellipses mark where the density is $1/e$ times that at the peak of the distribution. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

# Bayesian Discriminant Functions

*For Minimum Error:*

$$+ \ P(\omega_i | X)$$

$$+ \ P(X | \omega_i).P(\omega_i)$$

$$+ \ \log P(X | \omega_i) + \log P(\omega_i)$$

*For Minimum Risk:*

$$- \ R^i(X)$$

*Where*

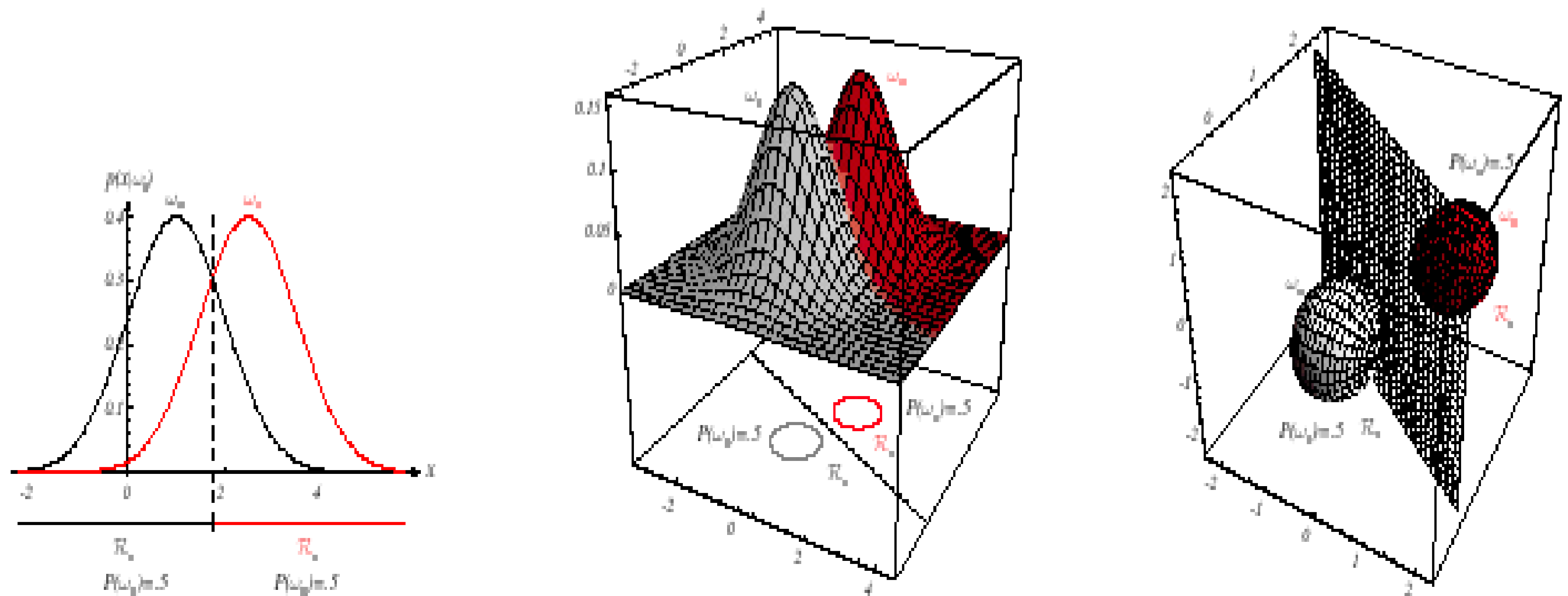$$R^i(X) = \sum_{j=1}^{c} \lambda(\alpha_i | \omega_j).P(\omega_j | X)$$

**FIGURE 2.10.** If the covariance matrices for two distributions are equal and proportional to the identity matrix, then the distributions are spherical in $d$ dimensions, and the boundary is a generalized hyperplane of $d-1$ dimensions, perpendicular to the line separating the means. In these one-, two-, and three-dimensional examples, we indicate $p(\mathbf{x}|\omega_i)$ and the boundaries for the case $P(\omega_1) = P(\omega_2)$. In the three-dimensional case, the grid plane separates $\mathcal{R}_1$ from $\mathcal{R}_2$. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

# *Bayes (Maximum Likelihood) Decision Classifier:*

Bayes Classifier can be shown to result with a minimum average error/risk, therefore considered to be optimal

*Most general optimal solution*

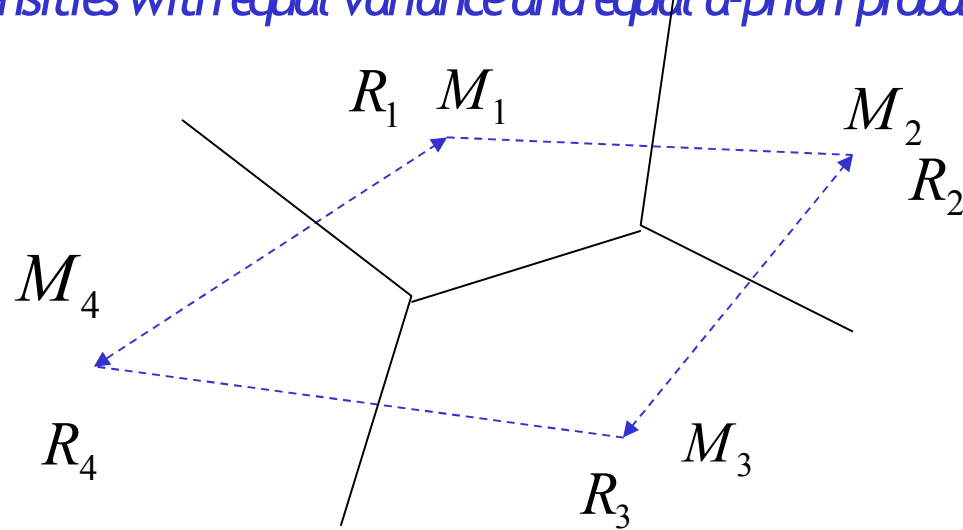*Can be used if the parametric models are known or properly estimated*

*Provides an upper limit(you cannot do better with other rule)*

*Useful in comparing with other classifiers*

# Minimum Distance Classifier: A special case of Bayes

*Classify an unknown sample X to the category with closest mean !*

*Optimum when gaussian densities with equal variance and equal a-priori probability.*



*Piecewise linear boundary in case of more than 2 categories.*

# Naive Bayes Classifier

Another special case

Assumes the features are statistically independent for a given category

Results with: $P(x_1, x_2, x_3...) = p(x_1)p(x_2)p(x_3)...$

Simplifies the decision rule

Often used in practice

# *Estimation of Parameters*

Bayes Rule is great if you know the class-conditional densities, but not available in nature

If the parametric form of the densities are given or assumed, then, using the labeled samples, the parameters can be estimated. (supervised learning)

Maximum Likelihood Estimation of parameters

Use the sample set $X_1, X_2, X_3, \ldots\ldots$

Find the parameters that will result with most likely combination

# *Features and Feature Selection*

**Curse of dimensionality:** High number of features increase the correct classification ratio but require too much data!

We should remove unnecessary features

Are all features independent from each other? Can we reduce the size without loosing information, by eliminating redundancies?(Principal Component analysis)
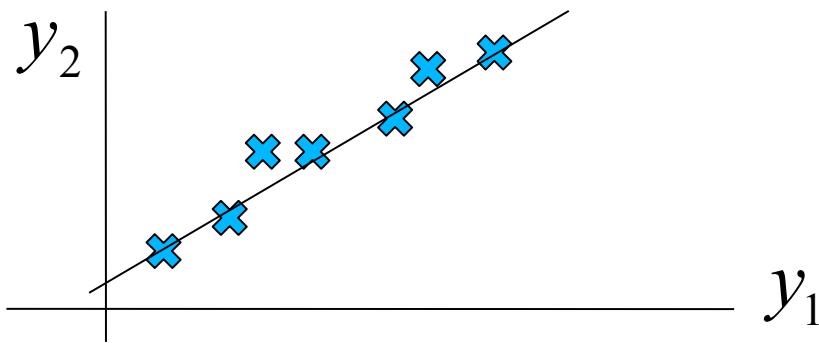
Are the features class discriminating? Again, remove features that have no discriminating ability(Fishers etc)

# Eliminating Redundant Features

$$X = \left[ x_1, \ldots \ldots, x_d \right]^T$$

is to be found using a larger set

$$Y = \left[ y_1, \ldots \ldots, y_m \right]^T$$



$$y_1, y_2$$

*Features that are linearly dependent so they can be reduced to 1*

*S*o we either

Throw one away

Generate a new feature using       and
  (ex:projections of the points to a line)

Form a linear combination of features.

$$x_1 = f_1(y_1, \ldots, y_m)$$
$$x_2 = f_2(y_1, \ldots, y_m)$$
$$x_d = f_d(y_1, \ldots, y_m)$$

*Linear functions of y*

$$X = WY$$

A linear transformation

W?  Can be found by:

**K-L expansion, Principal Component Analysis(PCA)**

PCA uses **eigenvalue** approach to result with an ordered set of
 features,
 in order of statistical independence.
So pick first d 'most independent' features

# Fisher's Linear Discriminant

- Curse of dimensionality. More features, more samples needed.
- We would like to choose features with more discriminating ability.
- Reduces the dimension of the  problem to one in its simplest form.
- Seperates samples from different categories.
- Consider samples from 2 different categories now.

Line 1(bad solution)

Blue (original data)
Red(projection onto Line 1)
Yellow(projection onto Line 2)

Line 2 (good solution)

-Find a line so that the projection separates
the samples best.

Same as: $\qquad y = W^T X$

Apply a transformation to samples X to result
with a scalar such that Fisher's criterion function is maximized, where

$$J(W) = \frac{(\mu_1 - \mu_2)^2}{\sigma_1^2 + \sigma_2^2}$$

It can be shown that W that maximizes J can be found
by solving the eigenvalue problem again:

$$S_W^{-1} S_B W = \lambda W$$

and the solution is given by

$$\boxed{W = S_W^{-1}(M_1 - M_2)}$$

Where SB and Sw are 'scatter matrices'
defined as functions of data scattering.

Multiple Discriminant Analysis:  c category problem.
A generalization of 2-category problem. Generalization to
M dimensions is also possible.

# k-Nearest Neighbor (k-NN) Rule

Non-parametric classification rules:
Linear and generalized discriminant functions
Nearest Neighbor & k-NN rules

Nearest Neighbor Classification  Rule
1-NN: A direct classification using learning samples
Assume we have learning samples from different categories

$$X^{11}, X^{12}, \ldots\ldots, X^{jk}, \ldots\ldots\ldots\ldots$$



$$X^{ik} \underline{\quad d \quad} X^{jl}$$

Assume a distance measure between samples such as euclidian

$$d(X^{ik}, X^{jl})$$

A general distance metric should obey the following rules:

$$d(X^{ij}, X^{ij}) = 0$$

$$d(X^{ij}, X^{jl}) = d(X^{jl}, X^{ij})$$

$$d(X, Y) \leq d(X, Z) + d(Z, Y)$$



Most standard: Euclidian Distance

$$d(X, Y) = \|X - Y\| = \left[ \sum_{i=1}^{n} (x_i - y_i)^2 \right]^{1/2} = \left[ (X - Y)^T (X - Y) \right]^{1/2}$$

1-NN Rule: Given an unknown sample X

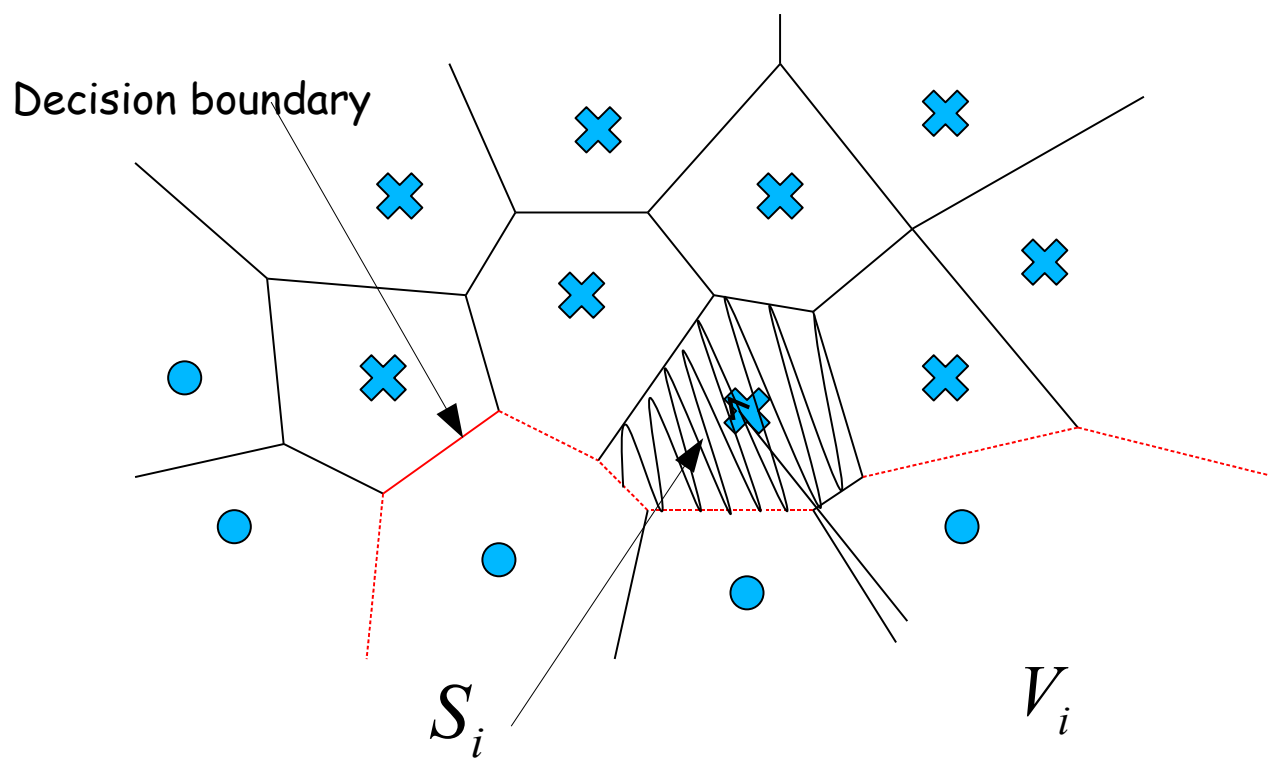$$\alpha_i \ \text{if}$$

$$d(X, X^{ik}) < d(X, X^{jl})$$

For $jl \neq ik$

That is, assign X to category $\omega_i$ if the closest neighbor of X is from category i.

Results with piecewise linear decision boundaries.

# Voronoi Diagrams

Decision boundary

$V_i$

$S_i$

$V_i$ is a polygon such that any point that falls in $V_i$ is closer to sample $S_i$ than any other sample $S_j$.

k-NN rule: instead of looking at the closest sample, we look at k nearest neighbors to X and we take a vote. The largest vote wins. k is usually taken as an odd number so that no ties occur.

Analysis of NN rule is possible when $M \to \infty$ and it was shown that it is no worse than twice of the <u>minimum-error classification (in error rate)</u>.

## EDITING AND CONDENSING

NN rule becomes very attractive because of its simplicity and yet good performance.

So, it becomes important to reduce the computational costs involved.

Do an intelligent elimination of the samples.



Remove samples that do not contribute to the decision boundary.

# *Linear Discriminant Functions*

Assume the discriminant functions are linear functions of X

$$g(X) = w_1 x_1 + w_2 x_2 + .... + w_n x_n + w_o$$
$$= W^T X + w_0$$
$$W = [w_1 ....... w_n]^T$$
$$X = [x_1 .......... x_n]^T$$

We may also write g in more compact form
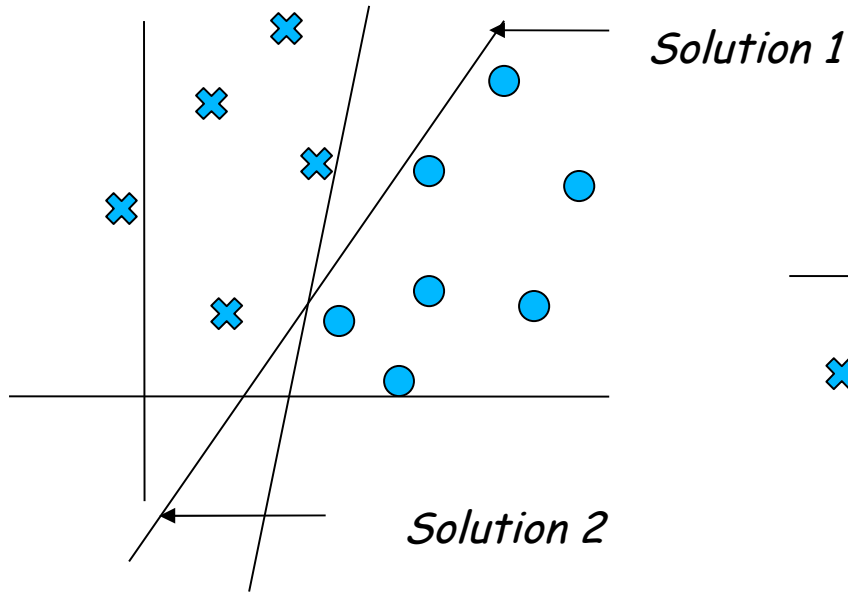
$$\boxed{g(X) = W_a^T X_a = W_a^T Y}$$

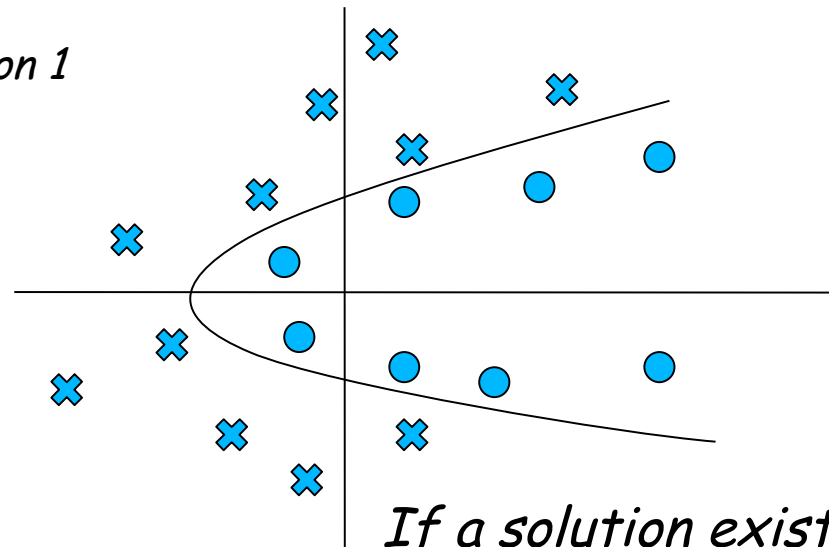$$Y = X_a = [x_1 ..... x_n 1]$$

$$W_a = [w_1 w_2 ...... w_n w_0]^T$$

# *Linear Discriminant Functions*

- Its assumed that the discriminant functions are linear (boundaries are linear)

- The labeled learning samples are used to find best boundaries.

- Finding the g is the same as finding W a.

- How do we define the 'best'?

- All learning samples are classified correctly?
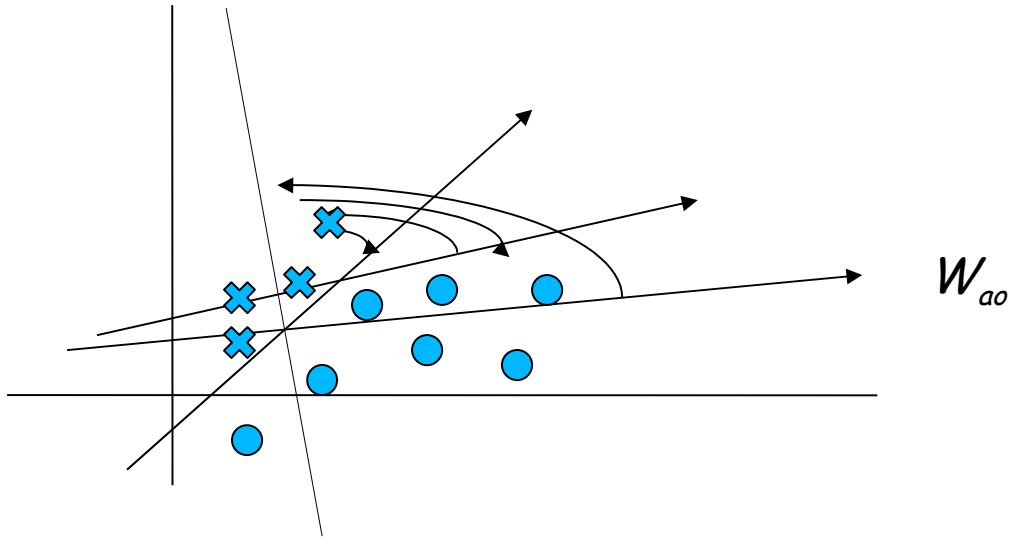
- Does a solution exist?

# Linear Separability

y

x

XOR Problem
Not linearly separable

Solution 1

Solution 2

Seperable,many solutionsions possible

If a solution exists-the problem is called "linearly separable" and $W_a$ is found iteratively. Otherwise "not linearly separable" piecewise or higher degree solutions are seeked.

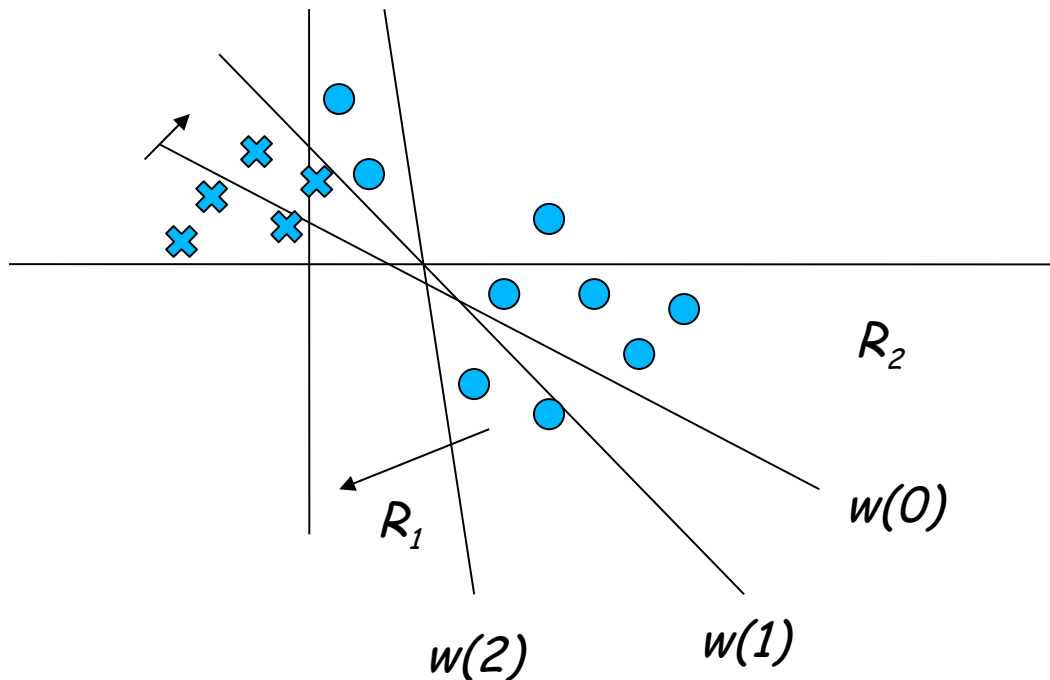# Iterative Solution: start with an initial estimate and update it until a solution is found.



$W_{ao}$

- **Gradient Descent Procedures**
- **Perceptron Criterion function**
- **One-layer Perceptron(Rosenblatt)**

Assume linear separability.

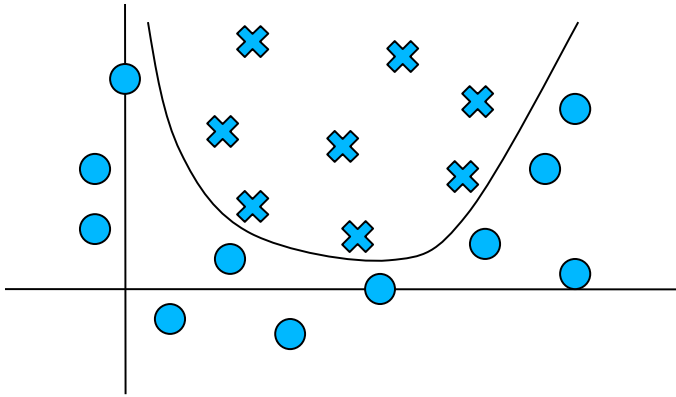<div align="center">

<span style="color:blue">perceptron learning</span>
</div>

An iterative algorithm that starts with an initial weight vector
and moves it back and forth until a solution is found.



Update of the weights are done using misclassified samples, towards
reducing them

# Generalized Discriminant Functions

When we have nonlinear problems as below:



Then we seek for a higher degree boundary.

Ex: quadratic boundary

$$g(X) = \sum \sum w_{ij} x_i x_j + \sum w_i x_i + w_0$$

will generate hyperquadratics boundaries.

g(X) still a linear function w's.
$$g(Y_a) = W_a^T Y_a$$

-

## NON-SEPARABLE CASE-What to do?

It was shown that we can increase the feature space dimension
with a nonlinear transformation, the results are linearly separable.
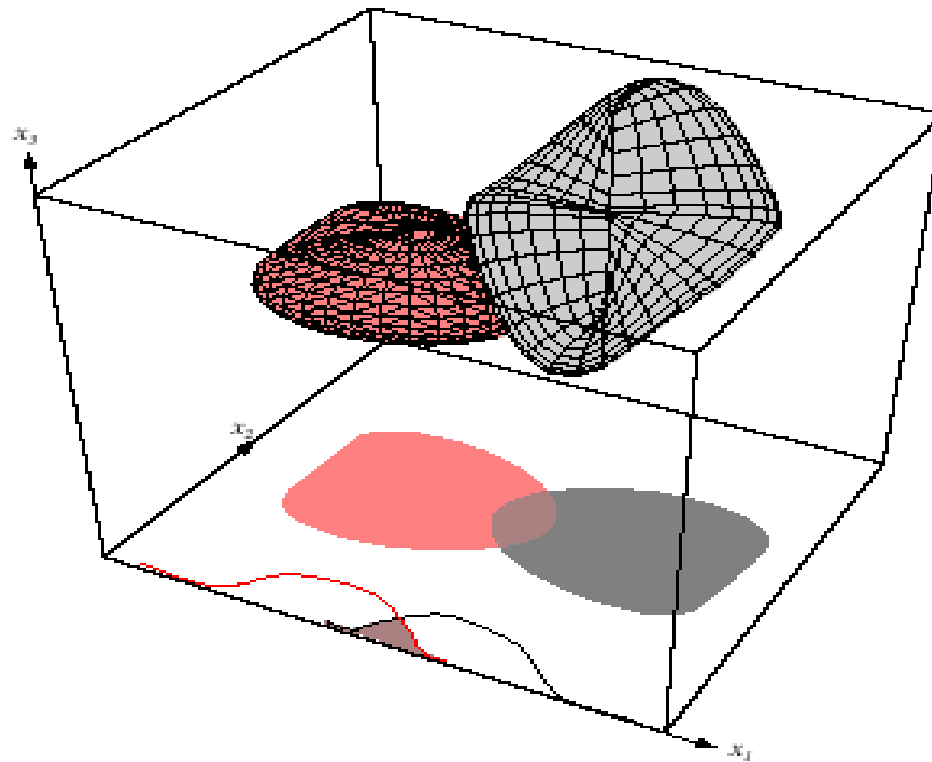Then find an optimum solution.(Support Vector Machines)

**FIGURE 3.3.** Two three-dimensional distributions have nonoverlapping densities, and thus in three dimensions the Bayes error vanishes. When projected to a subspace—here, the two-dimensional $x_1 - x_2$ subspace or a one-dimensional $x_1$ subspace—there can be greater overlap of the projected distributions, and hence greater Bayes error. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.
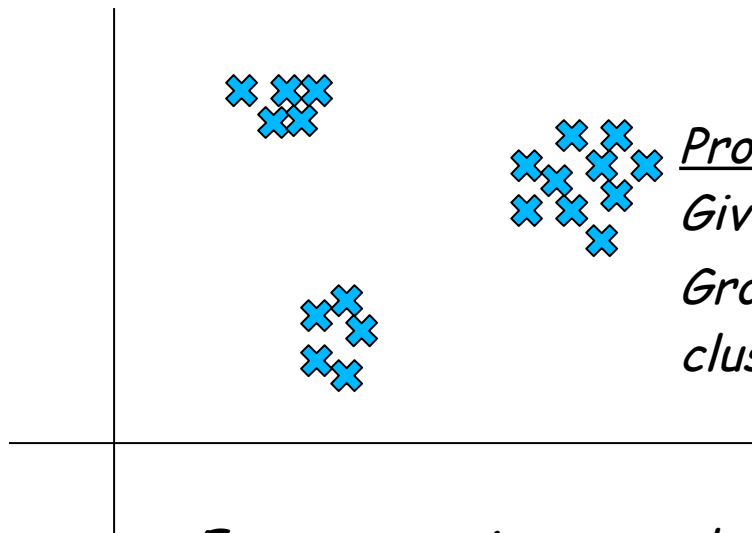
# UNSUPERVISED LEARNING AND CLUSTERING
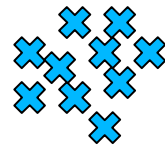
No class labels for learning samples.

We need additional means to label and classify – can be done separately (first label then classify) or together.

PARAMETRIC APPROACH- Estimation of class conditional densities
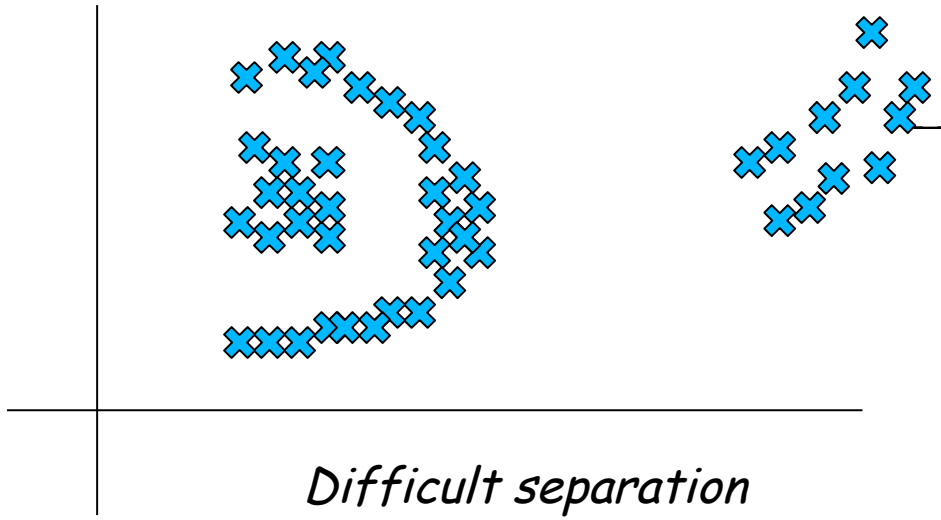
NONPARAMETRIC – CLUSTERING



*Problem:*

*Given samples $X_1, X_2, \ldots\ldots\ldots, X_n$*

*Group them into clusters so that samples in same cluster are "similar"*
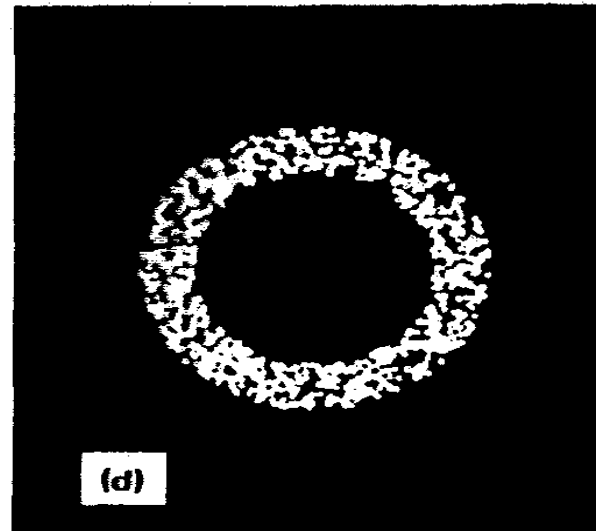
*Easy separation example*

*Difficult separation*

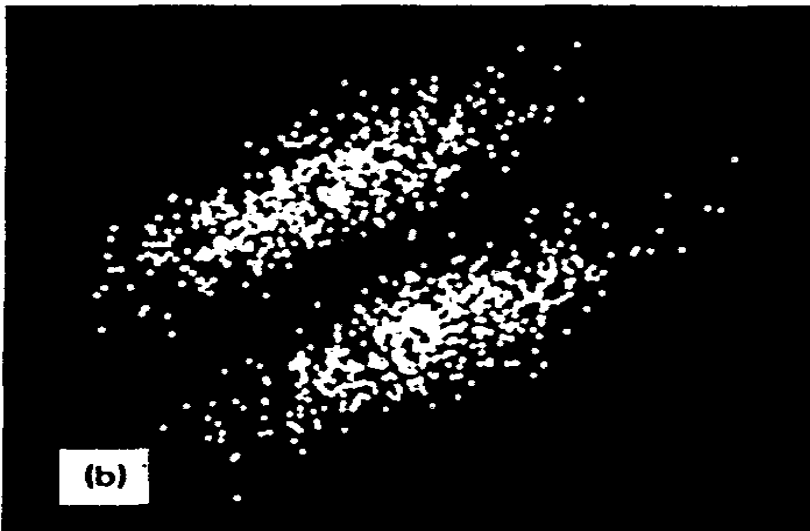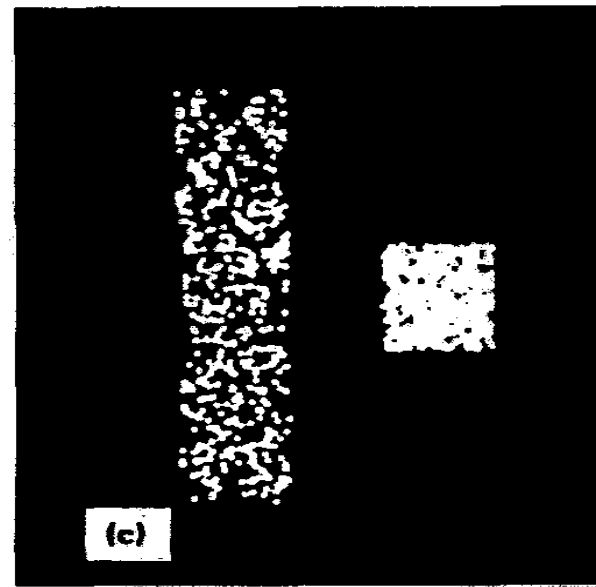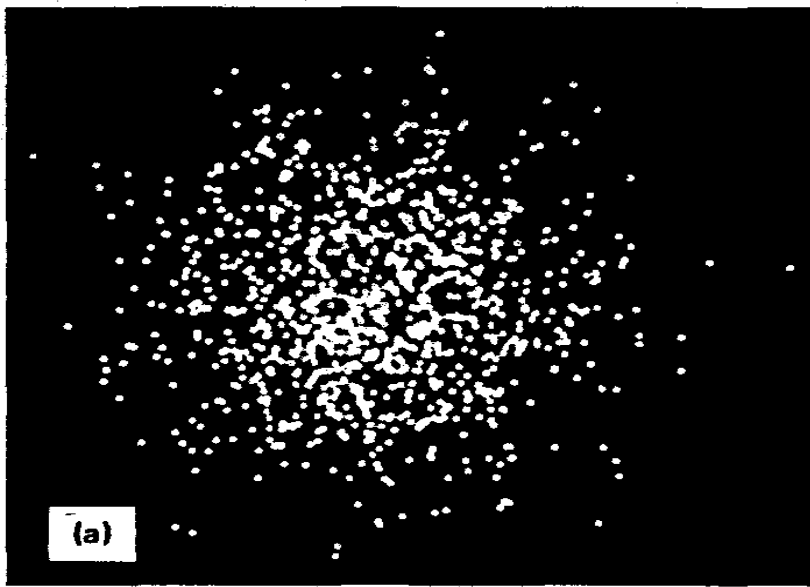**FIGURE 6.7.** Data sets having identical second-order statistics.
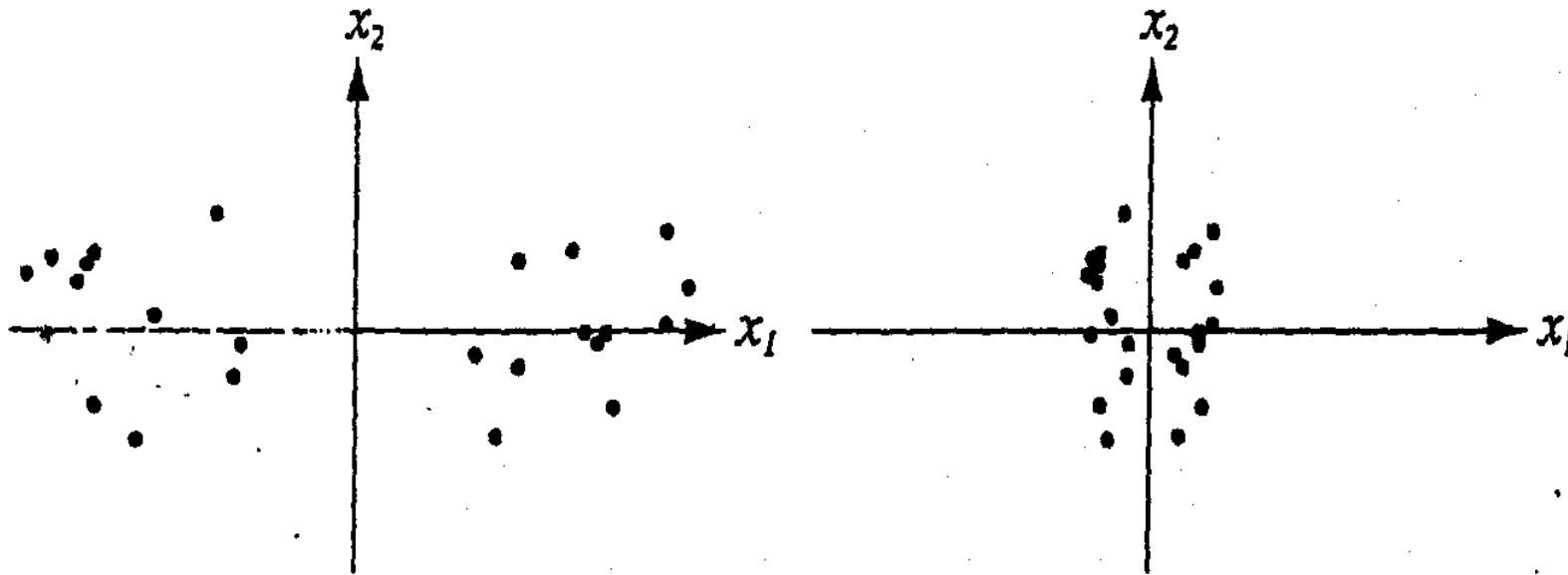
**FIGURE 10.9.** If the data fall into well-separated clusters (left), normalization by scaling for unit variance for the full data may reduce the separation, and hence be undesirable (right). Such a normalization may in fact be appropriate if the full data set arises from a single fundamental process (with noise), but inappropriate if there are several different processes, as shown here.

# Criterion Functions for Clustering

$$J_e = \sum_{i=1}^{c} \left( \sum_{x \in C_i} \|x - M_i\|^2 \right) = \frac{1}{2} \sum_{i=1}^{c} n_i \overline{S}_i$$

$$\overline{S}_i = \frac{1}{n_i^2} \sum_{x \in C_i} \sum_{x' \in C_i} \|x - x'\|^2$$

c- number of clusters

$M_i$ – mean of the samples in the same cluster

Aim: determine the partition that will minimize J.

Minimum variance partition ( sum of squared error criterion)

## Sum of squared error criterion

$$M_i = \frac{1}{n_i} \sum_{x \in D_i} x \qquad \textit{Sample mean for cluster } D_i$$

Sum of squared errors:

$$J_e = \sum_{i=1}^{c} \sum_{x \in D_i} \|x - M_i\|^2$$

Using $J_e$ results well for compact clusters.

## Basic ISODATA Algorithm (k-means)

Assume that there are k categories

- Choose k arbitrary points in space as cluster centers.

$$M_{10}, M_{20}, \ldots\ldots\ldots, M_{k0}$$

- Assign samples to their nearest cluster.

- Update M's. If any means changed value, go to 2. Otherwise stop.

May fall into local minimum.

---

■ **Algorithm 3.** (Basic Iterative Minimum-Squared-Error Clustering)

1 **begin initialize** $n, c, \mathbf{m}_1, \mathbf{m}_2, \ldots, \mathbf{m}_c$

2      **do** randomly select a sample $\hat{\mathbf{x}}$

3        $i \leftarrow \arg\min_{i'} \|\mathbf{m}_{i'} - \hat{\mathbf{x}}\|$      (classify $\hat{\mathbf{x}}$)

4        **if** $n_i \neq 1$ **then** compute

5       
$$\rho_j = \begin{cases} \frac{n_j}{n_j+1}\|\hat{\mathbf{x}} - \mathbf{m}_j\|^2 & j \neq i \\[2ex] \frac{n_j}{n_j-1}\|\hat{\mathbf{x}} - \mathbf{m}_i\|^2 & j = i \end{cases}$$

6        **if** $\rho_k \leq \rho_j$ for all $j$ **then** transfer $\hat{\mathbf{x}}$ to $\mathcal{D}_k$

7          recompute $J_e, \mathbf{m}_i, \mathbf{m}_k$

8      **until** no change in $J_e$ in $n$ attempts

9     **return** $\mathbf{m}_1, \mathbf{m}_2, \ldots, \mathbf{m}_c$

10 **end**

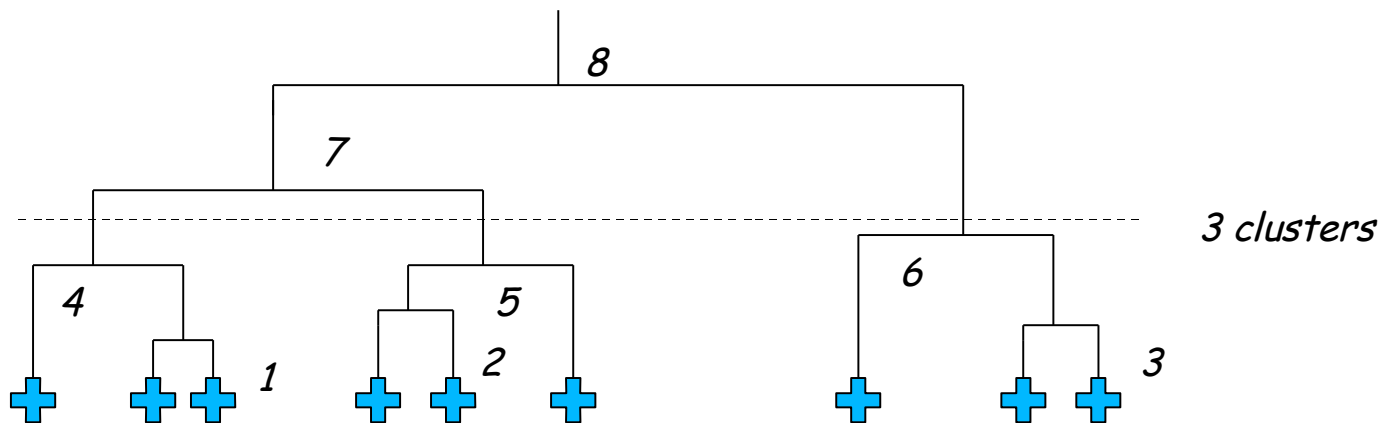---

54

# Hierarchical Clustering

A different approach to clustering.

Hierarchy in living species

Each species is a class by itself.

Combine the ones that are closest

Continue combining until the number of clusters are what is desired or a criterion is satisfied.



Issues:

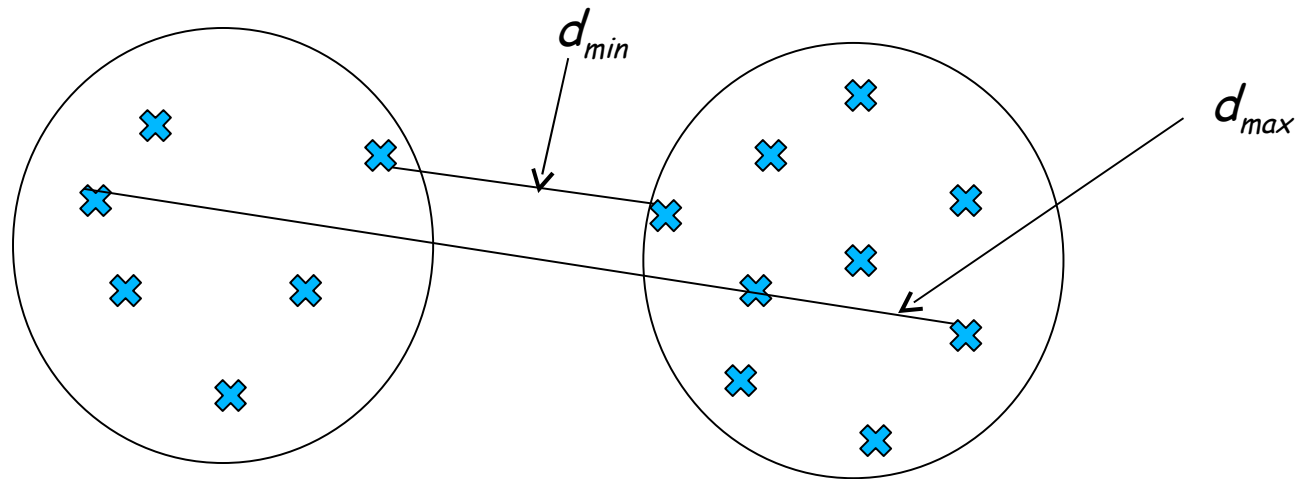How do we measure distance between clusters?

When do we stop?

Do we start from bottom or from top?

# Basic Bottom-up Hierarchical Clustering Algorithm for k clusters

**1.** Start taking each sample as a cluster. n=m (n of samples)

**2.** Measure $d_{i,j}$ – distance between clusters $D_i$ and $D_j$. Join two clusters $D_i$ and $D_k$ for which

n=n-1;

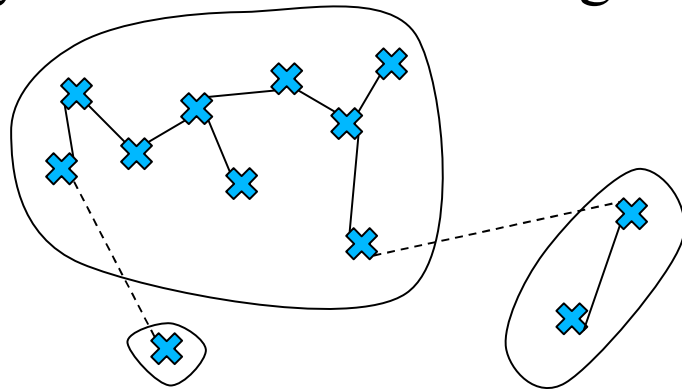**3.** If n<k   stop.          k:number of desired clusters
Else go to 2.

$d_{ij}$ can be defined in many ways.(where s is in i and z is in j:)

$$d_{\min_{i,j}} = \min \| x_s - x_z \|$$

$$d_{avg_{i,j}} = \frac{1}{n_i n_j} \left( \sum \sum \| x_s - x_z \| \right)$$

$$d_{\max} = \max \| x_s - x_z \|$$

$d_{min}$

$d_{max}$

## Example

Apply hierarchical clustering with $d_{min}$ to below data where c=3.



*will form elongated clusters!*

*Nearest Neighbor Clustering*

# Tree Classifiers

Consider the feature vector X = (x1, x2, x3....xn)

A tree classifier considers features one by one instead of as a whole and measures them one by one, following the leaves of a tree. The features are usually binary valued.

An optimum tree can be constructed using learning samples.

Leaves of the tree correspond to the classes.

Example will be seen in the following.

# Tree Classifiers: Example

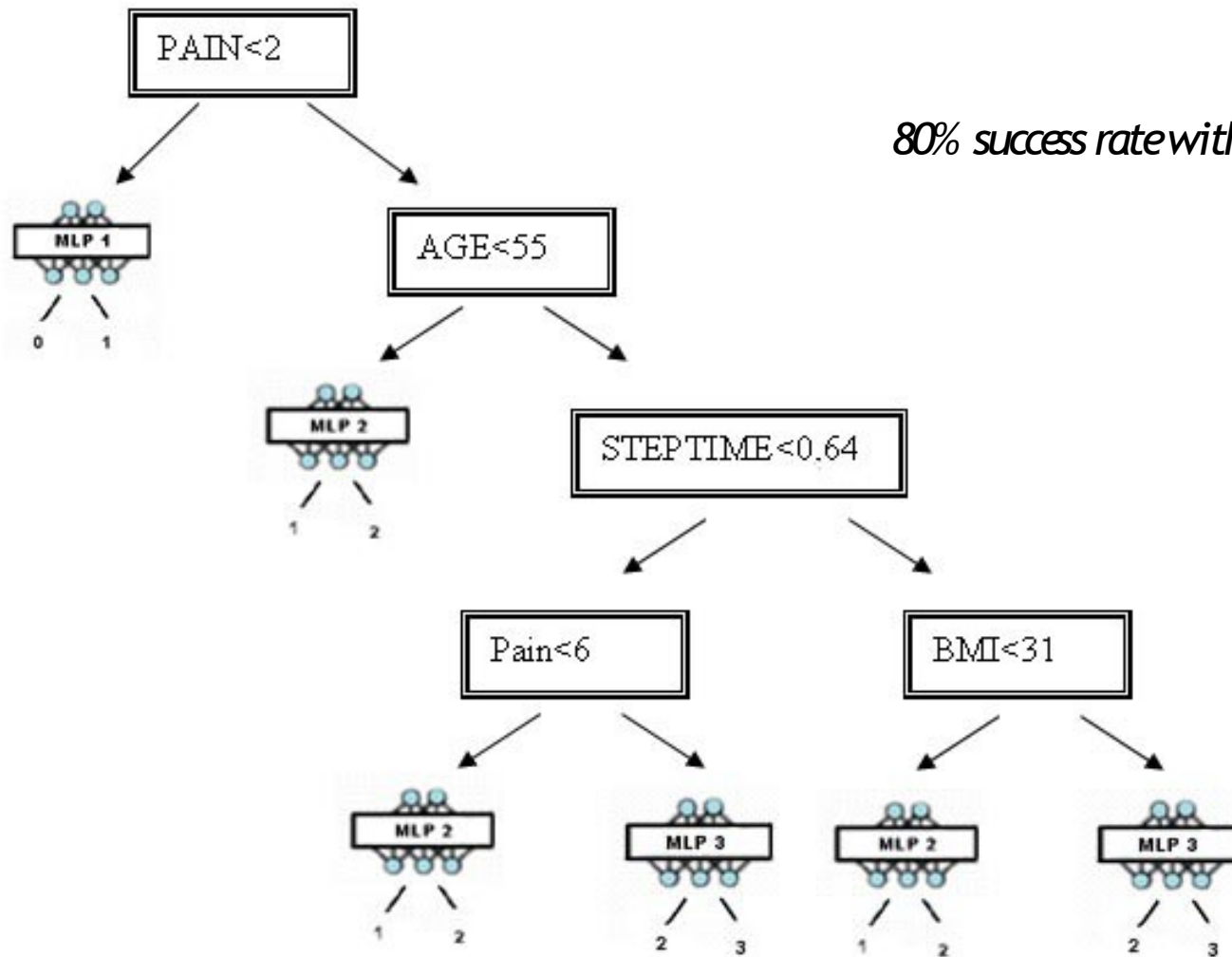Classification of illness 'Osteoarthritis' to levels 0–4

Using features obtained by gait (walking) analysis and patient history

The tree is generated by using 'gini' impurity index

*OAGAIT': A Decision Support System for Grading Knee Osteoarthritis using Gait Data' Pattern Recogniton Letters, July 2010*
*N. Köktaş, N. Yalabık, G. Yavuzer,P. Dunn, V. Atalay*

# Implementation and results



*80% success rate with 100 test samples*

# Decision Tree Construction as Learning

**Binary Tree** (All trees can be converted to a binary tree)

**'impurity' measure** used to decide how to split a tree (which feature to start with); shows how samples are distributed to categories as a result of split

**Split the tree so that impurity is lowest**

**'entropy impurity','variance impurity', 'Gini impurity'**

# *Training and Performance Testing in PR*

Whichever classifier is used, there is usually a training (learning) stage

How to train with available data?

Validation and Cross-Validation

How to test the performance?

Confusion Matrices

ROC Curves

# Cross-Validation

Assume a set of labeled samples are available

A number of them will be used to train the classifier and **others** to test the results (do not use the same samples)

**M -fold cross-validation:**

Divide the sample set into m disjoint sets of equal size

Train m times, each time with a different set used in testing

The performance is measured as the mean of errors each time

# *Performance Testing*

Simplest method: Confusion Matrices

Shows which category is confused with which

| | | | Estimated classes | | | | total | error rate |
|---|---|---|---|---|---|---|---|---|
| | | | 0 | 1 | 2 | 3 | total | error rate |
| Actual classes | 0 | | 27 | 7 | 5 | 1 | 40 | 0,325 |
| | 1 | | 3 | 24 | 7 | 6 | 40 | 0,4 |
| | 2 | | 3 | 8 | 21 | 8 | 40 | 0,475 |
| | 3 | | 3 | 5 | 12 | 20 | 40 | 0,5 |
| | Total | | 36 | 44 | 45 | 35 | 160 | 0,425 |

# ROC (Receiver Operating Characteristic) Curves

Assume 2 categories, where our aim is to detect a single object against all others (a binary classifier)

'hit' correctly classifying the object (true positive)

'false alarm' incorrectly finding that there is an object when it is not there (false positive)

'miss' finding no object when it is there (false negative)

Probability of 'hit' vs the probabilty of 'false alarm' is called a ROC Curve.

A ROC Curve is usually used to compare the classifiers.

|  | **P** | **N** |
|---|---|---|
| **Y** | True Positives | False Positives |
| **N** | False Negatives | True Negatives |

Hypothesized class

Column totals:   **P**   **N**

$$\text{fp rate} = \frac{FP}{N}$$

$$\text{tp rate} = \frac{TP}{P}$$

$$\text{precision} = \frac{TP}{TP+FP}$$

$$\text{recall} = \frac{TP}{P}$$

$$\text{accuracy} = \frac{TP+TN}{P+N}$$

$$\text{F-measure} = \frac{2}{1/\text{precision}+1/\text{recall}}$$

# Measures of Performance

## Sensitivity (recall rate)

$$\text{sensitivity} = \frac{\text{number of True Positives}}{\text{number of True Positives} + \text{number of False Negatives}}$$

## Specifity

$$\text{specificity} = \frac{\text{number of True Negatives}}{\text{number of True Negatives} + \text{number of False Positives}}$$

# Examples of ROC Curves

(Taken from :Tom Fawcett 'ROC Graphs: Notes and Practical
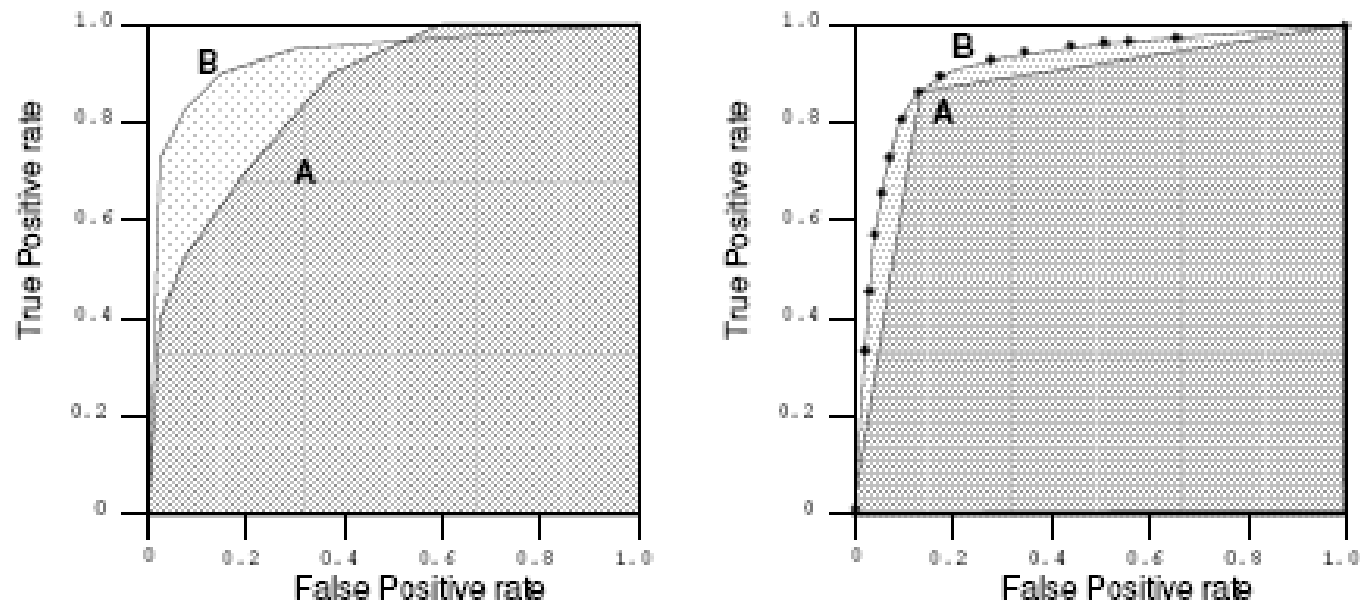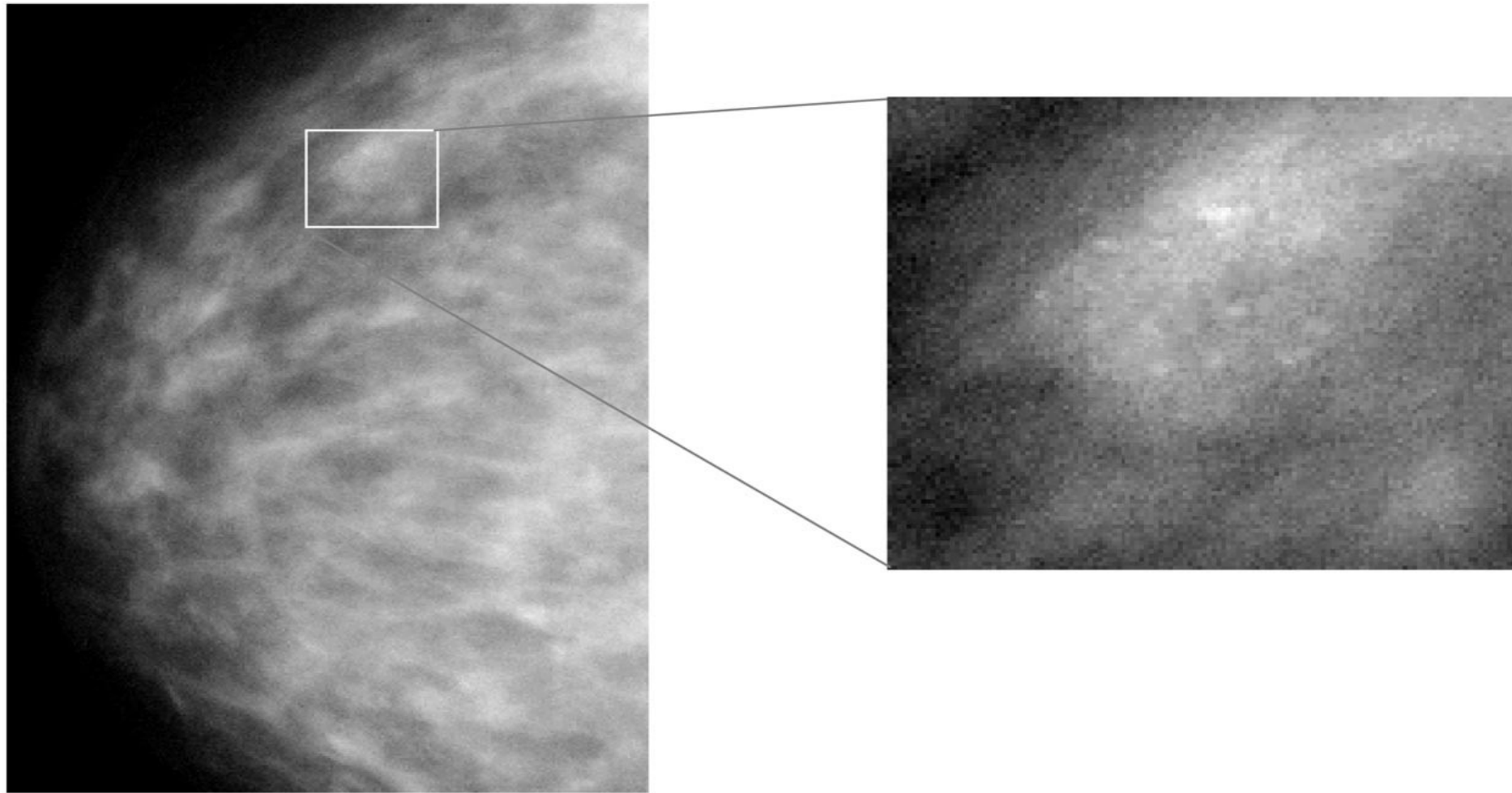Considerations for researchers')



Figure 7. Two ROC graphs. The graph on the left shows the area under two ROC curves. The graph on the right shows the area under the curves of a discrete classifier (A) and a probabilistic classifier (B).

# Example:Microcalcifications in a Mammogram

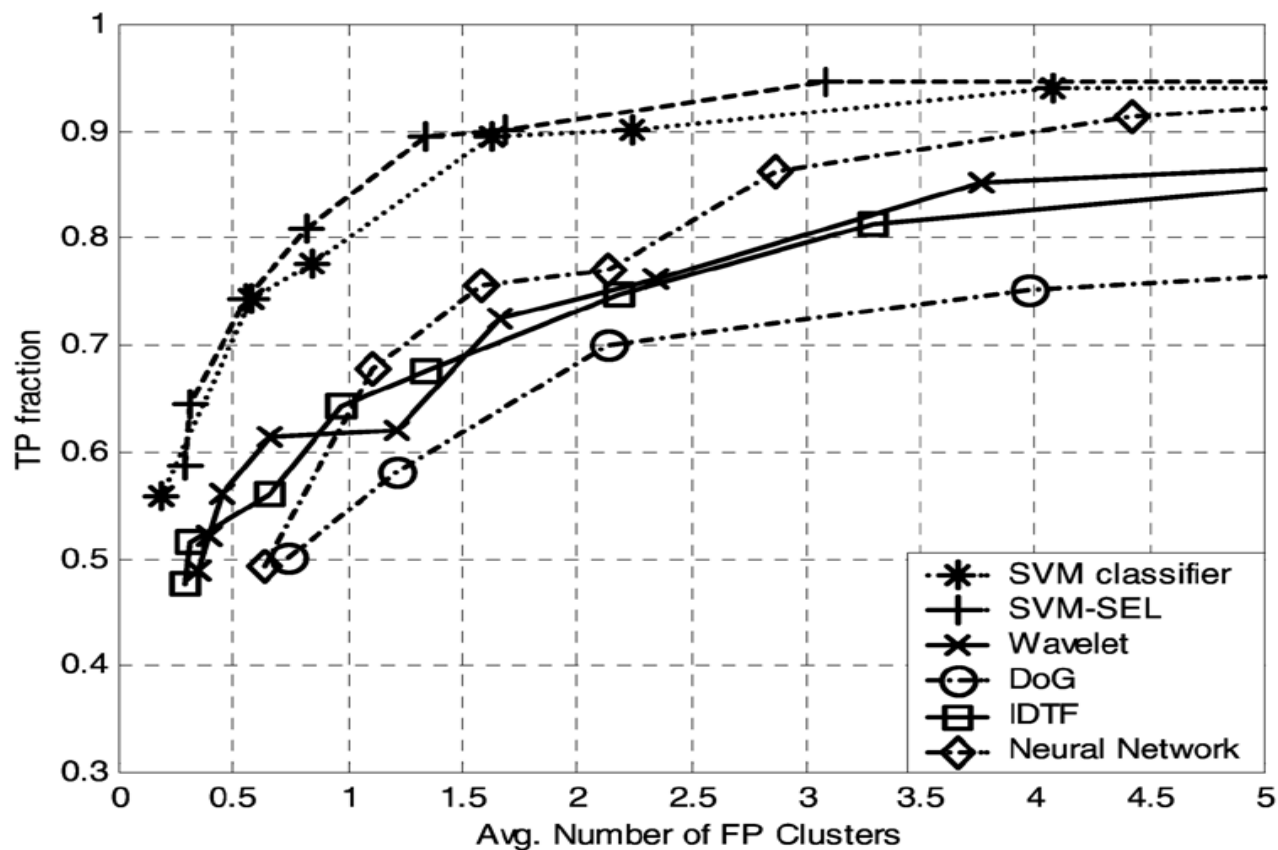'A Support Vector Machine Approach for Detection of Microcalcifications'

*Issam El-Naqa et al*

*IEEE TRANSACTIONS ON MEDICAL IMAGING, VOL. 21, NO. 12, DECEMBER 2002*

# Performance Comparison using a ROC curve

Higher the curve is, better the performance

# *References*

'*Pattern Classification*' R. Duda, P. Hart, D. Stork 2$^{nd}$ ed. Wiley 2001

*'Pattern Recognition'* S.Theodoridis, K.Koutroumbas" , Elsevier, 2003

Min720 Lecture Notes, N. Yalabık, ODTÜ 2010

http://home.comcast.net/~tom.fawcett/public_html/pap ROC Curves)

Others in respective pages