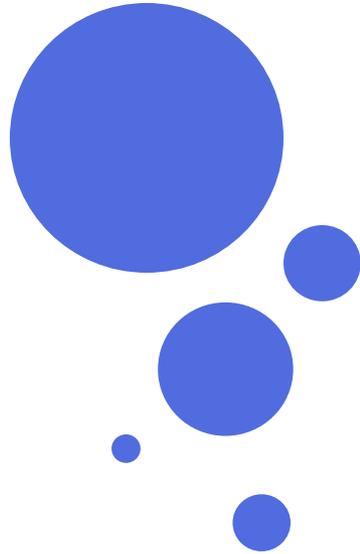




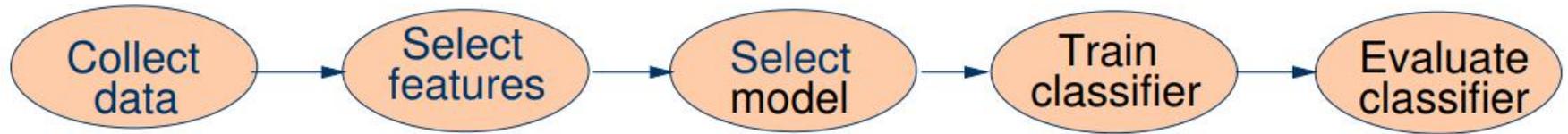
# Rensselaer

## Lecture 12: Midterm Exam Review

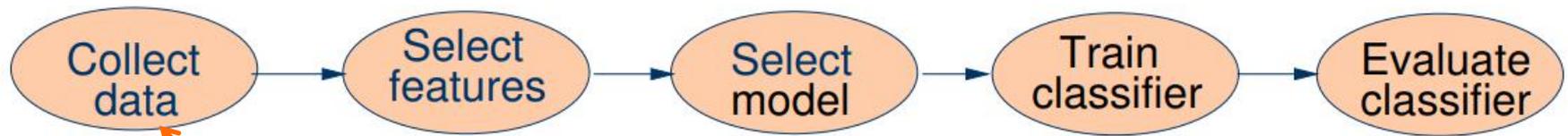


Dr. Chengjiang Long  
Computer Vision Researcher at Kitware Inc.  
Adjunct Professor at RPI.  
Email: [longc3@rpi.edu](mailto:longc3@rpi.edu)

# Pattern recognition design cycle



# Pattern recognition design cycle



- Collecting training and testing data.
- How can we know when we have adequately large and representative set of samples?

# Training/Test Split

- Randomly split dataset into two parts:
  - Training data
  - Test data
- Use training data to optimize parameters
- Evaluate error using test data

# Training/Test Split

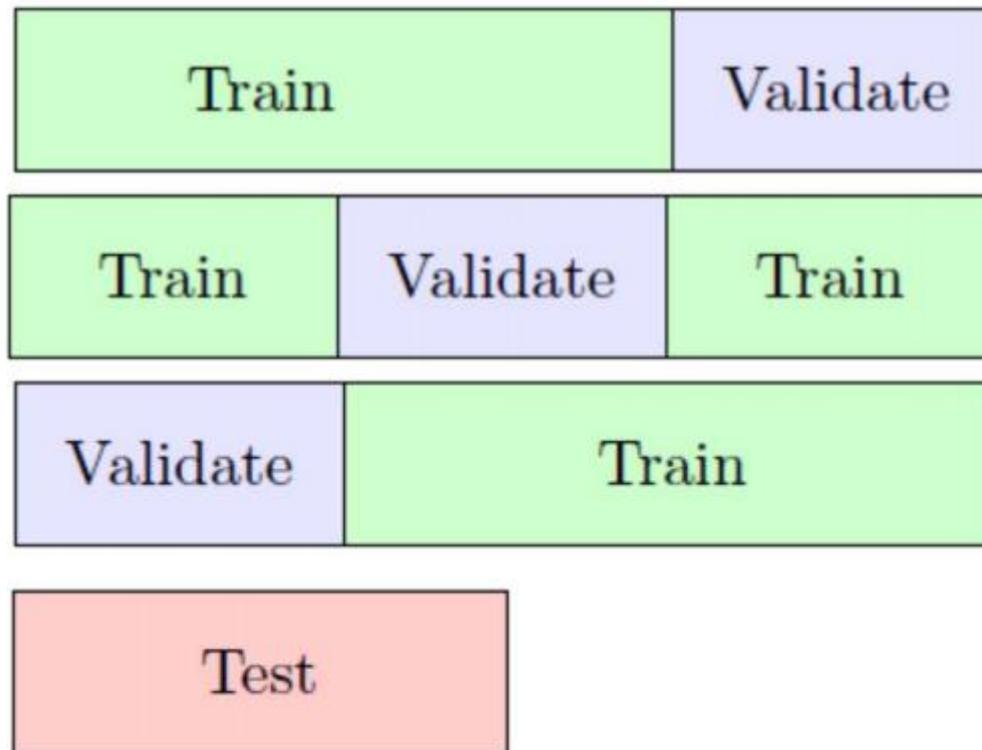
- How many points in each set ?
- Very hard question
  - Too few points in training set, learned classifier is bad
  - Too few points in test set, classifier evaluation is insufficient
- Cross-validation
- Leave-one-out cross-validation

# Cross-Validation

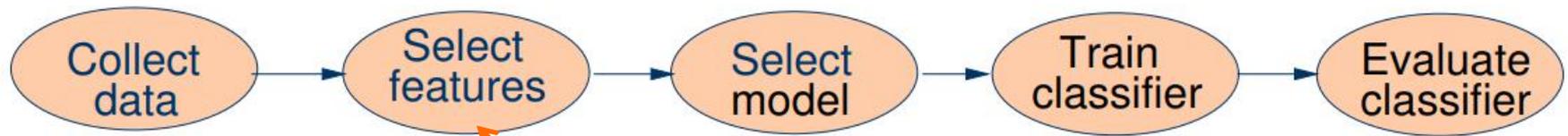
- In practice
- Available data  $\Rightarrow$  training and validation
- Train on the training data
- Test on the validation data
- k-fold cross validation:
  - Data randomly separated into k groups
  - Each time k
  - 1 groups used for training and one as testing

# Cross Validation and Test Accuracy

- Using CV on training + validation
- Classify test data with the best parameters from CV

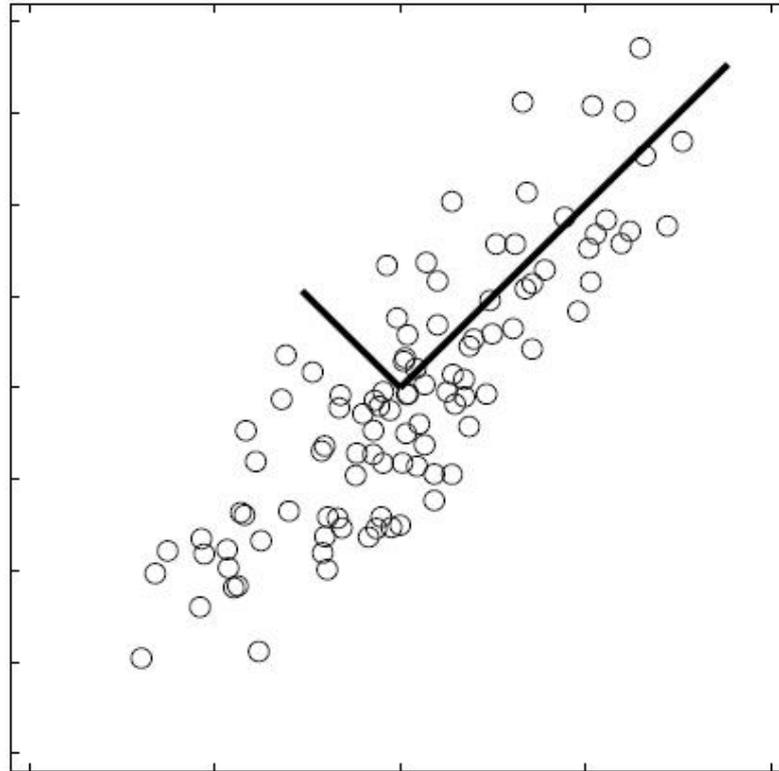


# Pattern recognition design cycle



- Domain dependence and prior information.
- Computational cost and feasibility.
- Discriminative features, i.e., similar values for similar patterns, and different values for different patterns.
- Invariant features with respect to translation, rotation and scale.
- Robust features with respect to occlusion, distortion, deformation, and variations in environment.

# PCA: Visualization



Data points are represented in a rotated **orthogonal** coordinate system: the origin is the **mean** of the data points and the axes are provided by the **eigenvectors**.

# Computation of PCA

- In practice we compute PCA via SVD (singular value decomposition)
- Form the centered data matrix:

$$X_{p,N} = [(\mathbf{x}_1 - \mathbf{m}) \dots (\mathbf{x}_N - \mathbf{m})]$$

- Compute its SVD:

$$X_{p,N} = U_{p,p} D_{p,p} (V_{N,p})^T$$

$U$  and  $V$  are orthogonal matrices,  $D$  is a diagonal matrix

# Computation of PCA...

- Sometimes we are given only a few high dimensional data points, *i.e.*,  $p \geq N$
- In such cases compute the SVD of  $X^T$ :

$$X^T = V_{N,N} D_{N,N} (U_{p,N})^T$$

So we get:

$$X = U_{p,N} D_{N,N} (V_{N,N})^T$$

Then, proceed as before, choose only  $d < N$  significant eigenvalues for data representation:

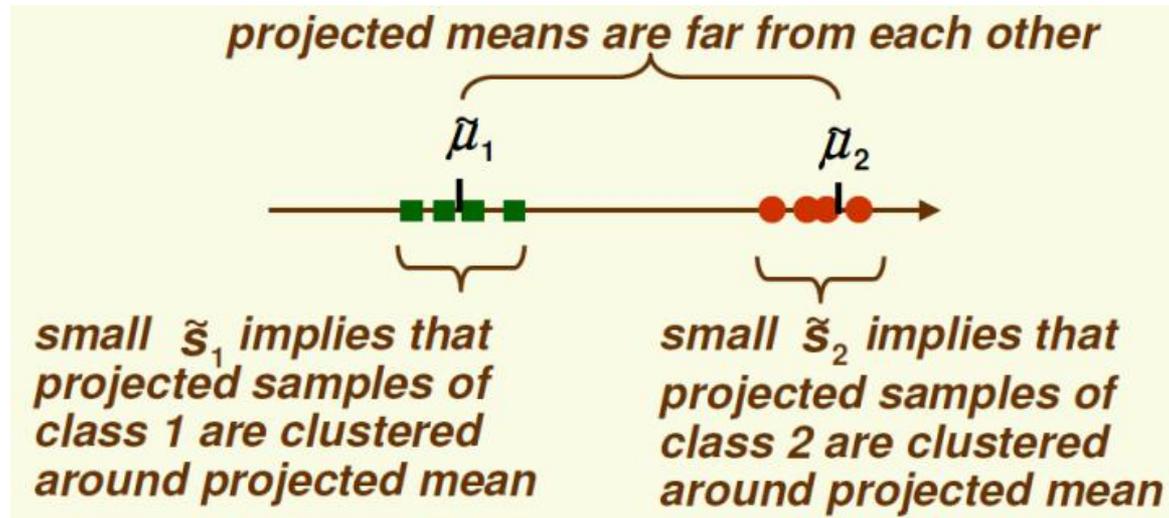
$$\tilde{\mathbf{x}}_i = \mathbf{m} + U_{p,d} (U_{p,d})^T (\mathbf{x}_i - \mathbf{m})$$

Usually we used the features with reduced dimensions to fit the classification models.

# Fisher Linear Discriminant

- We need to normalize by both scatter of class 1 and scatter of class 2
- The Fisher linear discriminant is the projection on a line in the direction  $\mathbf{v}$  which maximizes

$$J(\mathbf{v}) = \frac{(\tilde{\mu}_1 - \tilde{\mu}_2)^2}{\tilde{\mathbf{s}}_1^2 + \tilde{\mathbf{s}}_2^2}$$



# Fisher Linear Discriminant

- Thus our objective function can be written:

$$J(\mathbf{v}) = \frac{(\tilde{\mu}_1 - \tilde{\mu}_2)^2}{\tilde{\mathbf{s}}_1^2 + \tilde{\mathbf{s}}_2^2} = \frac{\mathbf{v}^t \mathbf{S}_B \mathbf{v}}{\mathbf{v}^t \mathbf{S}_W \mathbf{v}}$$

- Maximize  $J(\mathbf{v})$  by taking the derivative w.r.t.  $\mathbf{v}$  and setting it to 0

$$\begin{aligned} \frac{d}{d\mathbf{v}} J(\mathbf{v}) &= \frac{\left( \frac{d}{d\mathbf{v}} \mathbf{v}^t \mathbf{S}_B \mathbf{v} \right) \mathbf{v}^t \mathbf{S}_W \mathbf{v} - \left( \frac{d}{d\mathbf{v}} \mathbf{v}^t \mathbf{S}_W \mathbf{v} \right) \mathbf{v}^t \mathbf{S}_B \mathbf{v}}{(\mathbf{v}^t \mathbf{S}_W \mathbf{v})^2} \\ &= \frac{(2\mathbf{S}_B \mathbf{v}) \mathbf{v}^t \mathbf{S}_W \mathbf{v} - (2\mathbf{S}_W \mathbf{v}) \mathbf{v}^t \mathbf{S}_B \mathbf{v}}{(\mathbf{v}^t \mathbf{S}_W \mathbf{v})^2} = 0 \end{aligned}$$

# Fisher Linear Discriminant

Need to solve  $\mathbf{v}^t \mathbf{S}_W \mathbf{v} (\mathbf{S}_B \mathbf{v}) - \mathbf{v}^t \mathbf{S}_B \mathbf{v} (\mathbf{S}_W \mathbf{v}) = 0$

$$\Rightarrow \frac{\mathbf{v}^t \mathbf{S}_W \mathbf{v} (\mathbf{S}_B \mathbf{v})}{\mathbf{v}^t \mathbf{S}_W \mathbf{v}} - \frac{\mathbf{v}^t \mathbf{S}_B \mathbf{v} (\mathbf{S}_W \mathbf{v})}{\mathbf{v}^t \mathbf{S}_W \mathbf{v}} = 0$$

$$\Rightarrow \mathbf{S}_B \mathbf{v} - \frac{\mathbf{v}^t \mathbf{S}_B \mathbf{v} (\mathbf{S}_W \mathbf{v})}{\mathbf{v}^t \mathbf{S}_W \mathbf{v}} = 0$$

$$\Rightarrow \mathbf{S}_B \mathbf{v} = \lambda \mathbf{S}_W \mathbf{v}$$

*generalized eigenvalue problem*

# Fisher Linear Discriminant

$$\mathbf{S}_B \mathbf{v} = \lambda \mathbf{S}_W \mathbf{v}$$

- If  $\mathbf{S}_W$  has full rank (the inverse exists), we can convert this to a standard eigenvalue problem

$$\mathbf{S}_W^{-1} \mathbf{S}_B \mathbf{v} = \lambda \mathbf{v}$$

- But  $\mathbf{S}_B \mathbf{x}$  for any vector  $\mathbf{x}$ , points in the same direction as  $\mu_1 - \mu_2$

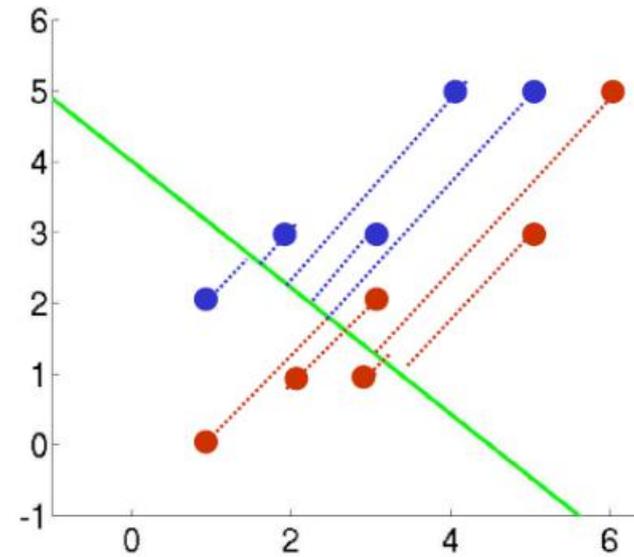
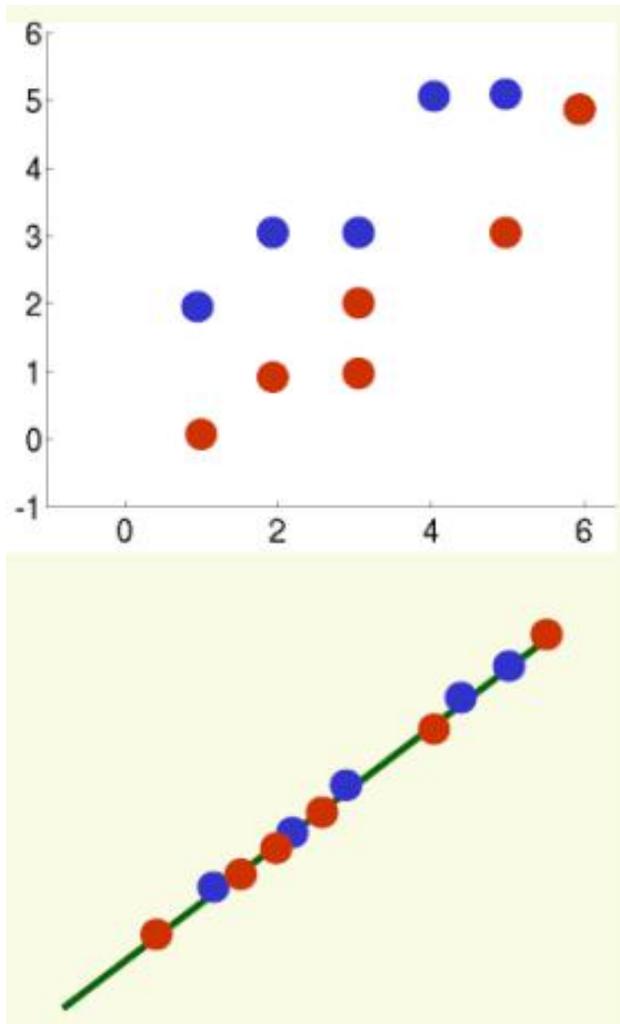
$$\mathbf{S}_B \mathbf{x} = (\mu_1 - \mu_2)(\mu_1 - \mu_2)^t \mathbf{x} = (\mu_1 - \mu_2) \underbrace{(\mu_1 - \mu_2)^t \mathbf{x}}_{\alpha} = \alpha (\mu_1 - \mu_2)$$

- Based on this, we can solve the eigenvalue problem directly

$$\mathbf{v} = \mathbf{S}_W^{-1} (\mu_1 - \mu_2)$$

$$\mathbf{S}_W^{-1} \mathbf{S}_B \underbrace{[\mathbf{S}_W^{-1} (\mu_1 - \mu_2)]}_{\mathbf{v}} = \mathbf{S}_W^{-1} [\alpha (\mu_1 - \mu_2)] = \underbrace{\alpha}_{\lambda} \underbrace{[\mathbf{S}_W^{-1} (\mu_1 - \mu_2)]}_{\mathbf{v}}$$

# Example



$$Y_1 = v^t c_1^t = [-0.65 \quad 0.73] \begin{bmatrix} 1 \cdots 5 \\ 2 \cdots 5 \end{bmatrix} = [0.81 \cdots 0.4]$$

$$Y_2 = v^t c_2^t = [-0.65 \quad 0.73] \begin{bmatrix} 1 \cdots 6 \\ 0 \cdots 5 \end{bmatrix} = [-0.65 \cdots -0.25]$$

# Pattern recognition design cycle



How can we know how close we are to the true model underlying the patterns?

- Domain dependence and prior information.
- Definition of design criteria.
- Parametric vs. non-parametric models.
- Handling of missing features.
- Computational complexity.
- Types of models: templates, decision-theoretic or statistical, syntactic or structural, neural, and hybrid.

# The Classifiers We Have Learned So Far

Bayesian classifiers {  
MLE classifier  
MAP classifier  
Naive Bayes classifier

Nonparametric classifiers → KNN classifier

Linear classifiers → LDF (Perceptron rule & Minimu Square Error rule & Ho-Kashyap Procedure) → SVM classifier

Nonlinear classifiers  $\xrightarrow[\text{Feature Mapping } \Phi]{\text{Kernel Tricks}}$  Linear classifiers

# Decision Rule

- Using **Bayes' rule**:

$$P(\omega_j / x) = \frac{p(x / \omega_j)P(\omega_j)}{p(x)} = \frac{\text{likelihood} \times \text{prior}}{\text{evidence}}$$

where

$$p(x) = \sum_{j=1}^2 p(x / \omega_j)P(\omega_j)$$

**Decide**  $\omega_1$  if  $P(\omega_1 / x) > P(\omega_2 / x)$ ; otherwise **decide**  $\omega_2$

or

**Decide**  $\omega_1$  if  $p(x/\omega_1)P(\omega_1) > p(x/\omega_2)P(\omega_2)$ ; otherwise **decide**  $\omega_2$

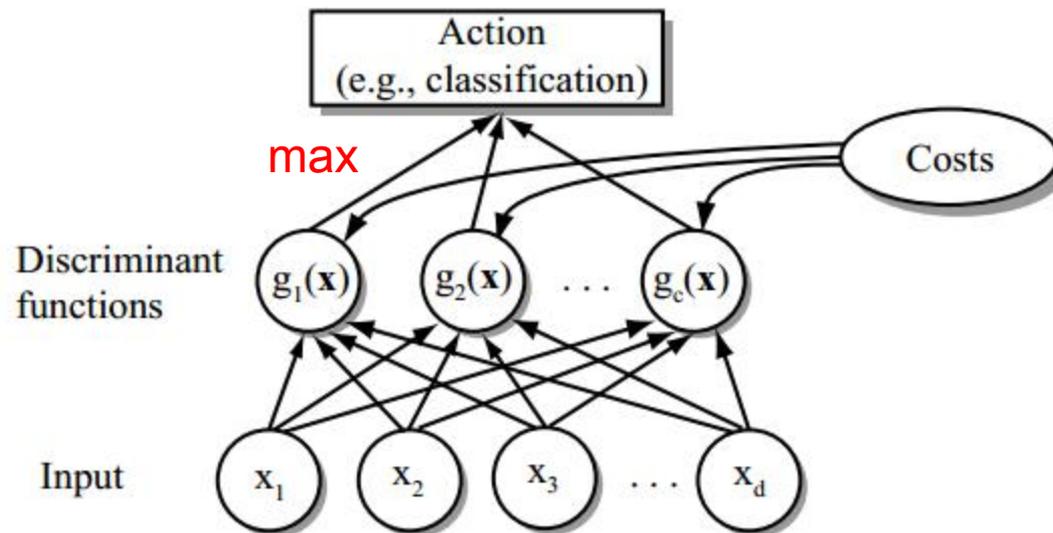
or

**Decide**  $\omega_1$  if  $p(x/\omega_1)/p(x/\omega_2) > P(\omega_2)/P(\omega_1)$  ; otherwise **decide**  $\omega_2$

# Discriminant Functions

- A useful way to represent a classifier is through **discriminant functions**  $g_i(\mathbf{x})$ ,  $i = 1, \dots, c$ , where a feature vector  $\mathbf{x}$  is assigned to class  $\omega_i$  if

$$g_i(\mathbf{x}) > g_j(\mathbf{x}) \text{ for all } j \neq i$$



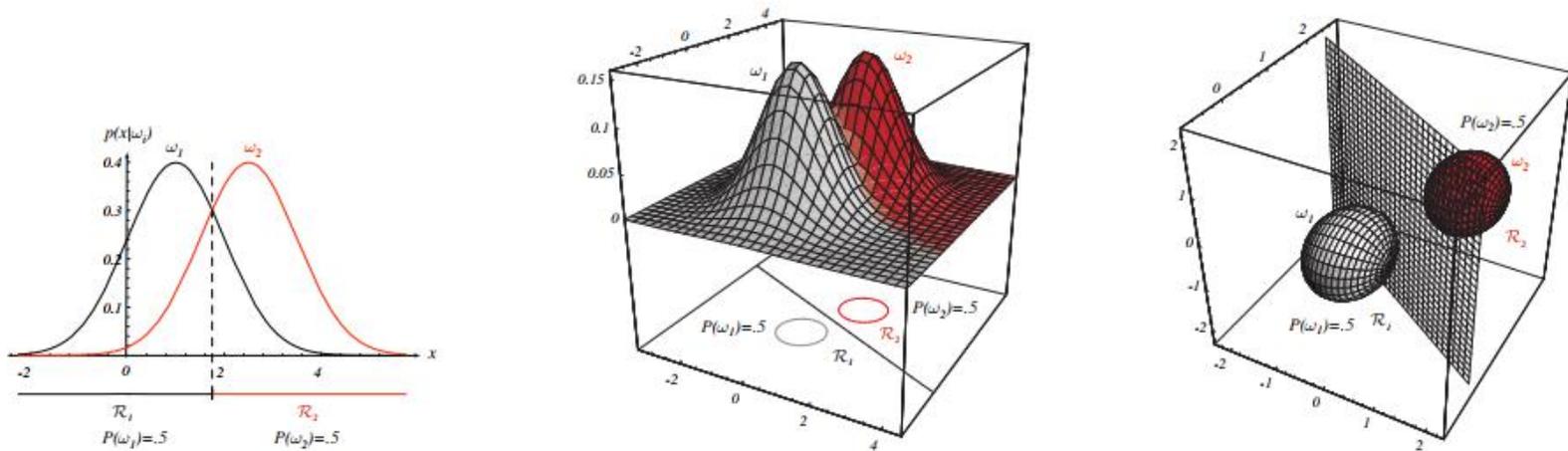
# Discriminants for Bayes Classifier

- Is the choice of  $g_i$  unique?
  - Replacing  $g_i(\mathbf{x})$  with  $f(g_i(\mathbf{x}))$ , where  $f(\cdot)$  is **monotonically increasing**, does not change the classification results.

$$g_i(\mathbf{x}) = P(\omega_i / \mathbf{x}) \quad \Rightarrow \quad g_i(\mathbf{x}) = \frac{p(\mathbf{x} / \omega_i) P(\omega_i)}{p(\mathbf{x})}$$
$$g_i(\mathbf{x}) = p(\mathbf{x} / \omega_i) P(\omega_i)$$
$$g_i(\mathbf{x}) = \ln p(\mathbf{x} / \omega_i) + \ln P(\omega_i)$$

we'll use this  
discriminant extensively!

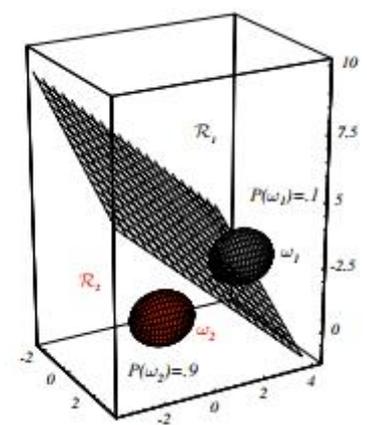
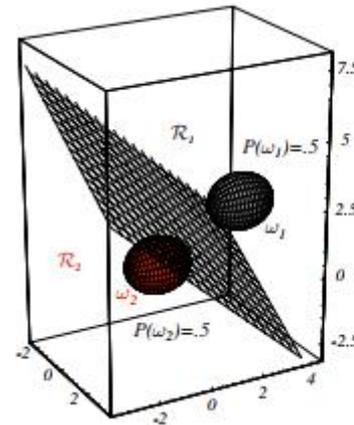
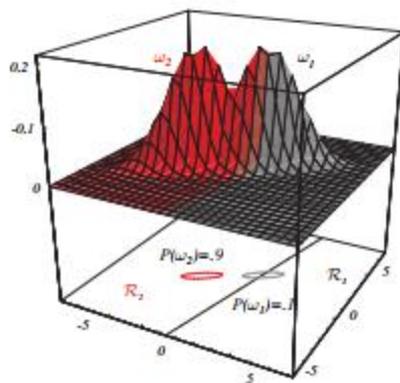
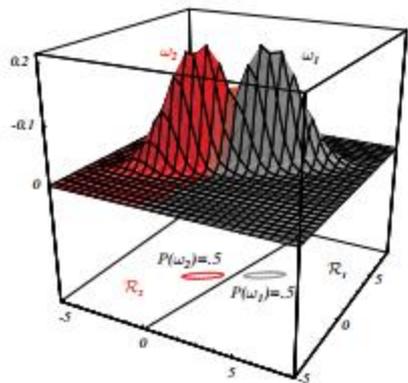
# Case 1: Statistically Independent Features with Identical Variances



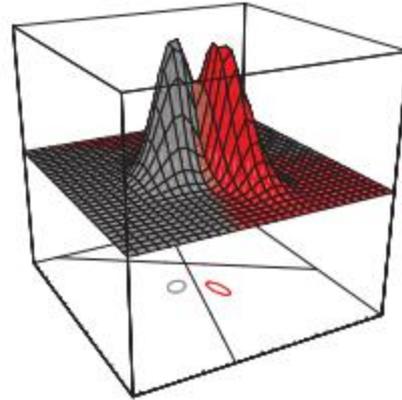
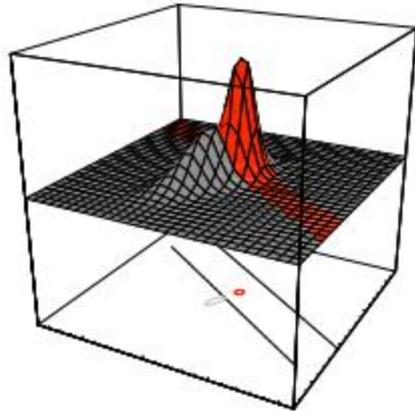
**FIGURE 2.10.** If the covariance matrices for two distributions are equal and proportional to the identity matrix, then the distributions are spherical in  $d$  dimensions, and the boundary is a generalized hyperplane of  $d - 1$  dimensions, perpendicular to the line separating the means. In these 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.

# Case II: Identical Covariances $\Sigma_i = \Sigma$

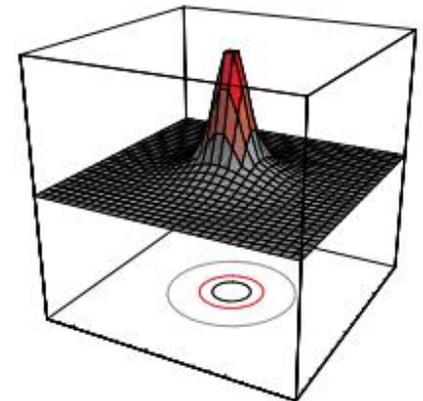
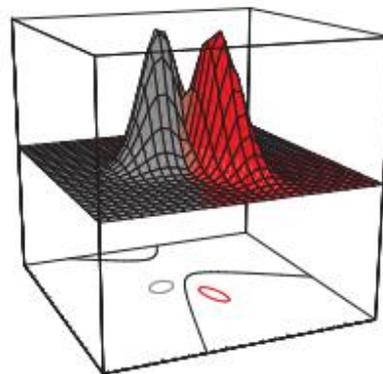
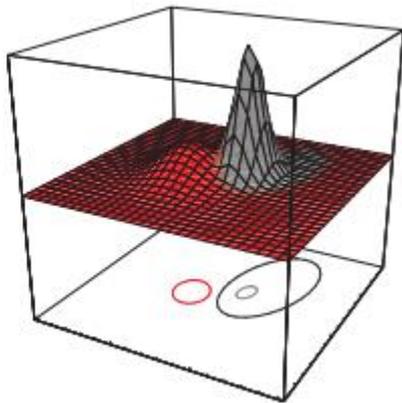
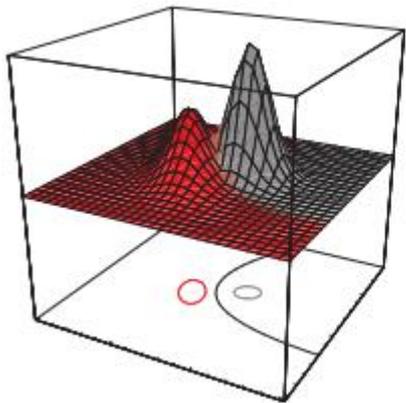
- Notes on Decision Boundary
  - As for Case I, passes through point  $x_0$  lying on the line between the two class means. Again,  $x_0$  in the middle if priors identical
  - Hyperplane defined by boundary generally not orthogonal to the line between the two means



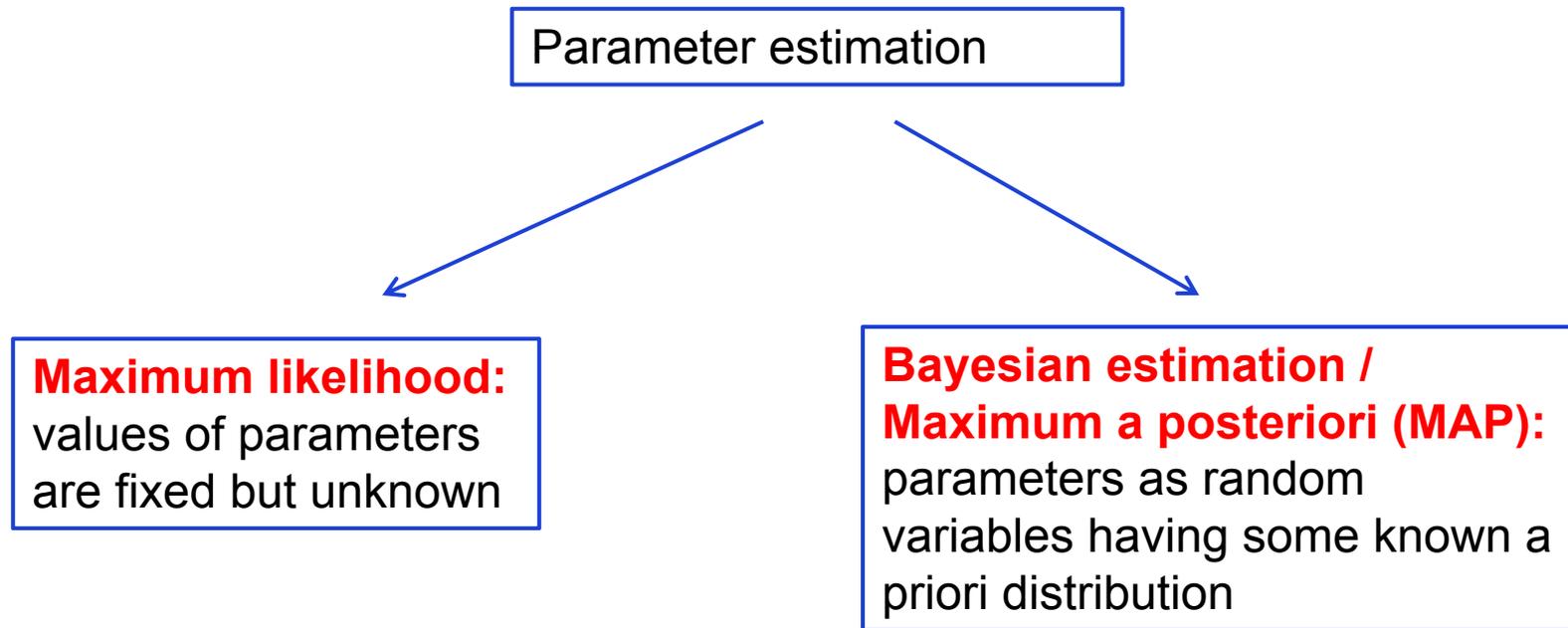
# Case III: arbitrary $\Sigma_i$



Nonlinear decision boundaries



# Parameter Parameter



# Maximum-Likelihood Estimation

- Use set of independent samples to estimate  $p(D | \theta)$

$$\text{Let } D = \{x_1, x_2, \dots, x_n\} \quad |D| = n$$

$$p(D | \theta) = \prod_{k=1}^{k=n} p(x_k | \theta)$$

- Our goal is to determine  $\hat{\theta}$  (value of  $\theta$  that best agrees with observed training data)
- Note if  $D$  is fixed  $p(D | \theta)$  is not a density

# Example: Gaussian case

- Assume we have  $c$  classes and

$$p(\mathbf{x} | \omega_j) \sim N(\mu_j, \Sigma_j)$$

$$p(\mathbf{x} | \omega_j) \equiv p(\mathbf{x} | \omega_j, \theta_j)$$

- Use the information provided by the training samples to estimate  $\theta = (\theta_1, \theta_2, \dots, \theta_c)$ , each  $\theta_j$  is associated with each category.
- Suppose that  $D$  contains  $n$  samples,  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$

$$p(D | \theta) = \prod_{k=1}^{k=n} p(\mathbf{x}_k | \theta)$$

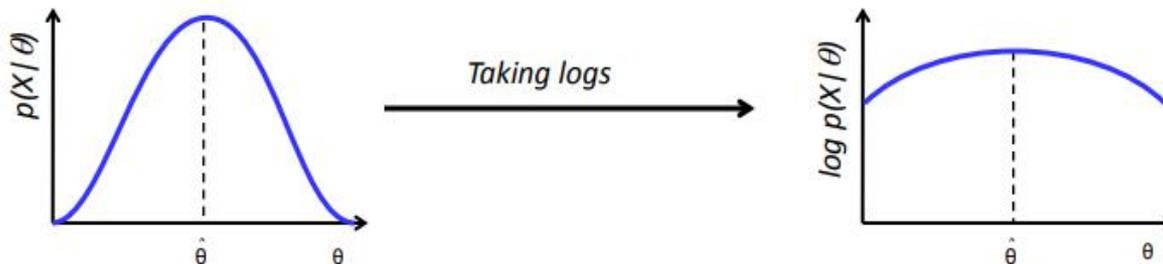
# Maximum-Likelihood Estimation

$$p(D | \theta) = \prod_{k=1}^{k=n} p(x_k | \theta)$$

- $p(D | \theta)$  is called the likelihood of  $\theta$  w.r.t the set of samples.
- ML estimate of  $\theta$  is, by definition the value  $\hat{\theta}$  that maximizes  $p(D | \theta)$

**“It is the value of  $\theta$  that best agrees with the actually observed training sample”**

$$\hat{\theta} = \operatorname{argmax}[p(X|\theta)] = \operatorname{argmax}[\log p(X|\theta)]$$



# Optimal Estimation

- Let  $\theta = (\theta_1, \theta_2, \dots, \theta_p)^t$  and let  $\nabla_{\theta}$  be the gradient operator

$$\nabla_{\theta} = \left[ \frac{\partial}{\partial \theta_1}, \frac{\partial}{\partial \theta_2}, \dots, \frac{\partial}{\partial \theta_p} \right]^t$$

- We define  $l(\theta)$  as the log likelihood function

$$l(\theta) = \ln p(D | \theta)$$

- New problem statement:
  - determine  $\theta$  that maximizes the log likelihood

$$\hat{\theta} = \arg \max_{\theta} l(\theta)$$

# Optimal Estimation

$$\nabla_{\theta} l = \sum_{k=1}^{k=n} \nabla_{\theta} \ln p(x_k | \theta)$$

$$\nabla_{\theta} l = 0$$

- Local or global maximum
- Local or global minimum
- Saddle point
- Boundary of parameter space

# Bayesian Estimation (MAP): General Theory

- $p(x | D)$  computation can be applied to any situation in which the unknown density can be parameterized.
- The basic assumptions are:
  - ✓ The form of  $p(x | \theta)$  assumed known, but the value of  $\theta$  is not known exactly
  - ✓ Our knowledge about  $\theta$  is assumed to be contained in a known prior density
  - ✓ The rest of our knowledge  $\theta$  is contained in a set  $D$  of  $n$  random variables  $x_1, x_2, \dots, x_n$  that follows  $p(x)$

# Bayesian Estimation (MAP): General Theory

- The basic problem is: **“Compute the posterior density  $p(\theta | D)$ ” then “Derive  $p(x | D)$  ”**

$$p(x | D) = \int \overset{\text{known}}{p(x | \theta)} \overset{\text{unknown}}{p(\theta | D)} d\theta$$

- Using Bayes formula, we have:

$$p(\theta | D) = \frac{p(D | \theta)p(\theta)}{\int p(D | \theta)p(\theta)d\theta}$$

- And by the independence assumption:

$$p(D | \theta) = \prod_{k=1}^{k=n} p(x_k | \theta)$$

# MLE vs. MAP

- Maximum Likelihood estimation (MLE)
  - Choose value that maximizes the probability of observed data

$$\hat{\theta}_{MLE} = \arg \max_{\theta} P(D|\theta)$$

- Maximum a posteriori (MAP) estimation
  - Choose value that is most probable given observed data and prior belief

$$\begin{aligned}\hat{\theta}_{MAP} &= \arg \max_{\theta} P(\theta|D) \\ &= \arg \max_{\theta} P(D|\theta)P(\theta)\end{aligned}$$

**When is MAP same as MLE?**

# Naïve Bayes Classifier (not BE)

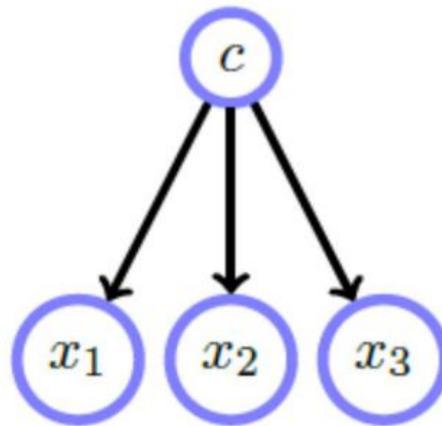
- Simple classifier that applies Bayes' rule with strong (naive) independence assumptions
- A.k.a. the 'independent feature model'

$$p(\omega_i | \mathbf{x}_1, \mathbf{x}_2, \dots) = \alpha p(\mathbf{x}_1 | \omega_i) p(\mathbf{x}_2 | \omega_i) \dots p(\omega_i)$$

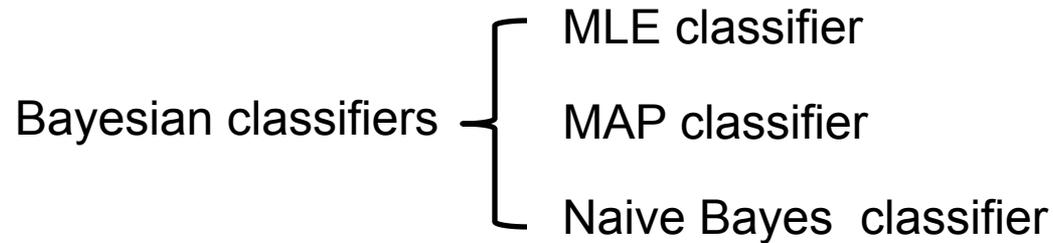
- Often performs reasonably well despite simplicity

# Naïve Bayes Classifier

- NB is known to produce posteriors closer to extremes (0 or 1) than true posteriors
  - Why?
- NB performs well when only small amounts of training data are available
  - Why?



# The Classifiers We Have Learned So Far



Nonparametric classifiers → KNN classifier

Linear classifiers → LDF (Perceptron rule & Minimu Square Error rule & Ho-Kashyap Procedure) → SVM classifier

Nonlinear classifiers  $\xrightarrow[\text{Feature Mapping } \Phi]{\text{Kernel Tricks}}$  Linear classifiers

# Density Estimation: Two Approaches

$$p(\mathbf{x}) \approx \frac{k/n}{V}$$

- Parzen Windows:

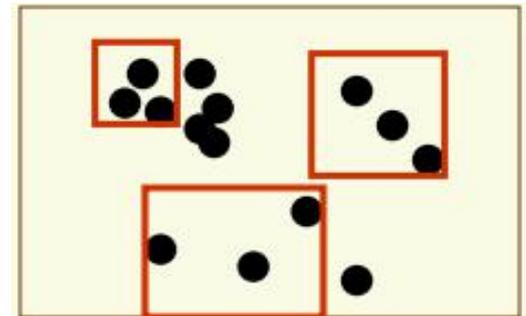
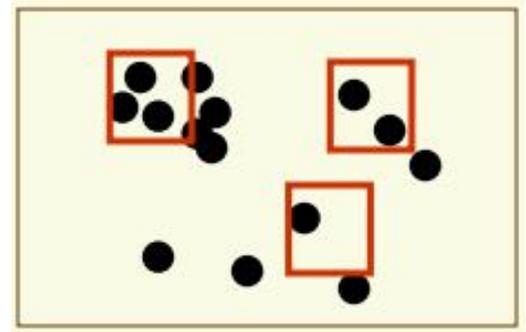
- ✓ Shrink an initial region where  $V_n = 1/\sqrt{n}$  and show that

$$p_n(\mathbf{x}) \xrightarrow{n \rightarrow \infty} p(\mathbf{x})$$

- ✓ This is called “the Parzen window estimation method”

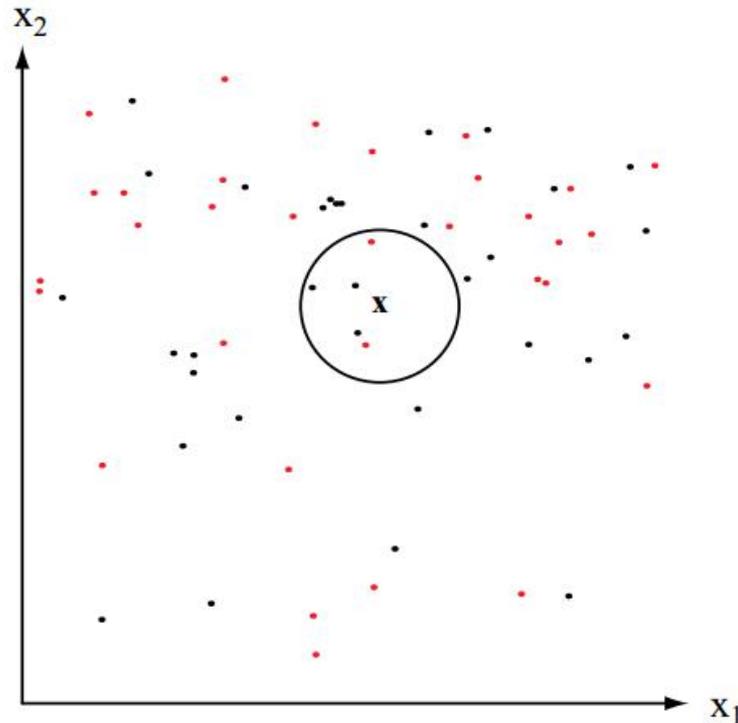
- k-Nearest Neighbors

- ✓ Specify  $k_n$  as some function of  $n$ , such as  $k_n = \sqrt{n}$  the volume  $V_n$  is grown until it encloses  $k_n$  neighbors of  $\mathbf{x}$ . This is called “the  $k_n$ -nearest neighbor estimation method”



# The k-Nearest-Neighbor Rule

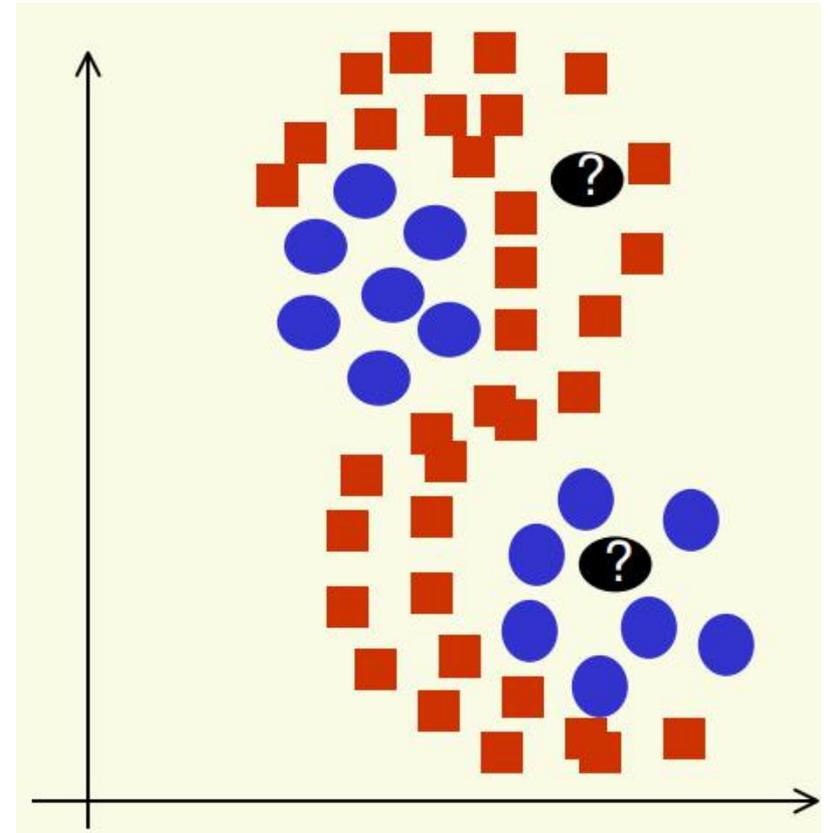
- **Goal:** Classify  $x$  by assigning it the label most frequently represented among the  $k$  nearest samples
- Use a voting scheme



The k-nearest-neighbor query starts at the test point and grows a spherical region until it encloses  $k$  training samples, and labels the test point by a majority vote of these samples

# kNN: Multi-modal Distributions

- Most parametric distributions would not work for this 2 class classification problem
- Nearest neighbors will do reasonably well, provided we have a lot of samples



# The Classifiers We Have Learned So Far

Bayesian classifiers { MLE classifier  
MAP classifier  
Naive Bayes classifier

Nonparametric classifiers → KNN classifier

Linear classifiers → LDF (Perceptron rule & Minimu Square Error rule & Ho-Kashyap Procedure) → SVM classifier

Nonlinear classifiers  $\xrightarrow[\text{Feature Mapping } \Phi]{\text{Kernel Tricks}}$  Linear classifiers

# Augmented feature space

- **Augmented** feature/parameter space

$$g(\mathbf{x}) = \mathbf{w}^t \mathbf{x} + w_0 = \sum_{i=1}^d w_i x_i + x_0 w_0 = \sum_{i=0}^d w_i x_i = \boldsymbol{\alpha}^t \mathbf{y}$$

$$\mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \\ \dots \\ w_d \end{bmatrix} \Rightarrow \boldsymbol{\alpha} = \begin{bmatrix} w_0 \\ w_1 \\ \dots \\ w_d \end{bmatrix} \quad \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_d \end{bmatrix} \Rightarrow \mathbf{y} = \begin{bmatrix} x_0 \\ x_1 \\ \dots \\ x_d \end{bmatrix}$$

**Discriminant:**

$$g(\mathbf{x}) = \boldsymbol{\alpha}^t \mathbf{y}$$

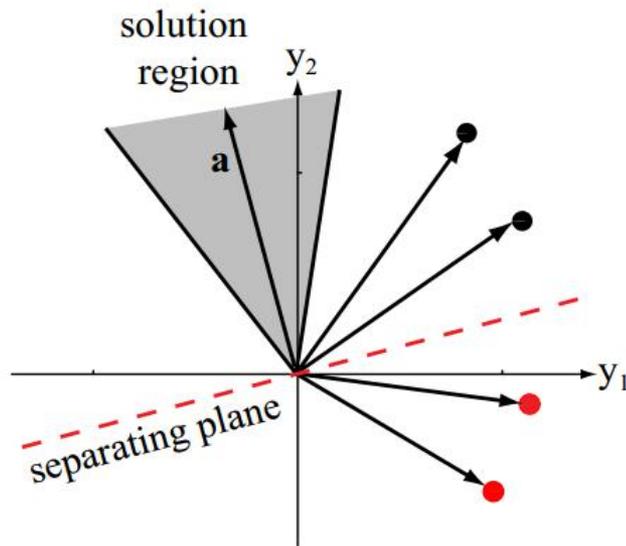
If	$\boldsymbol{\alpha}^t \mathbf{y}_i \geq 0$	assign $\mathbf{y}_i$ to $\omega_1$
else if	$\boldsymbol{\alpha}^t \mathbf{y}_i < 0$	assign $\mathbf{y}_i$ to $\omega_2$

# Normalization

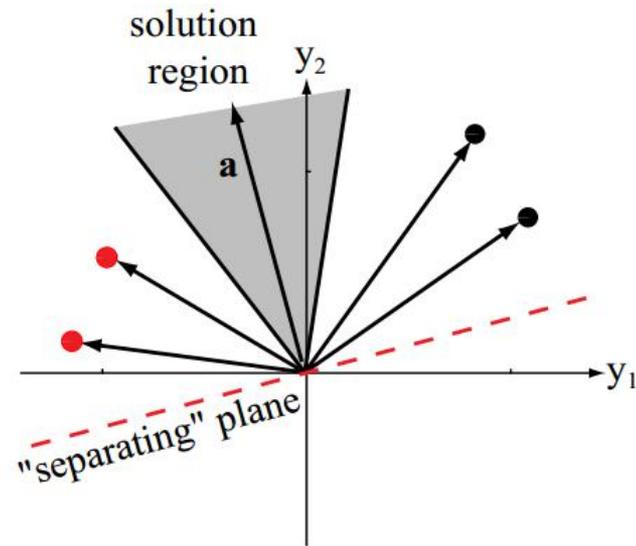
## Classification rule:

If  $\alpha^t \mathbf{y}_i > 0$  assign  $\mathbf{y}_i$  to  $\omega_1$   
else if  $\alpha^t \mathbf{y}_i < 0$  assign  $\mathbf{y}_i$  to  $\omega_2$

- If  $\mathbf{y}_i$  in  $\omega_2$ , replace  $\mathbf{y}_i$  by  $-\mathbf{y}_i$
- Find  $\alpha$  such that:  $\alpha^t \mathbf{y}_i > 0$



Seek a hyperplane that **separates** patterns from different categories



Seek a hyperplane that puts normalized patterns on the **same (positive) side**

# Perceptron Batch Rule

- The gradient of  $J_p(\mathbf{\alpha})$  is:

$$J_p(\mathbf{\alpha}) = \sum_{\mathbf{y} \in Y(\mathbf{\alpha})} (-\mathbf{\alpha}^t \mathbf{y}) \quad \Rightarrow \quad \nabla J_p = \sum_{\mathbf{y} \in Y(\mathbf{\alpha})} (-\mathbf{y})$$

- It is not possible to solve analytically  $\nabla J_p = 0$ .
- The perceptron update rule is obtained using gradient descent:

$$\mathbf{\alpha}(k+1) = \mathbf{\alpha}(k) + \eta(k) \sum_{\mathbf{y} \in Y(\mathbf{\alpha})} \mathbf{y}$$

- It is called batch rule because it is based on all misclassified examples

# Perceptron Single Sample Rule

- The gradient decent single sample rule for  $J_p(\mathbf{a})$  is:

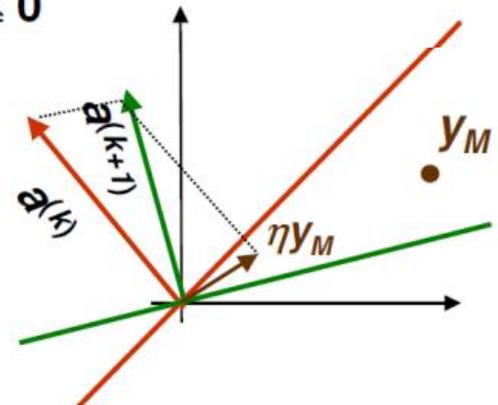
$$\mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} + \eta^{(k)} \mathbf{y}_M$$

- Note that  $y_M$  is one sample misclassified by  $\mathbf{a}^{(k)}$
- Must have a consistent way of visiting samples

- Geometric Interpretation:

- Note that  $y_M$  is one sample misclassified by  $(\mathbf{a}^{(k)})^t \mathbf{y}_M \leq 0$
- $y_M$  is on the wrong side of decision hyperplane
- Adding  $\eta y_M$  to  $\mathbf{a}$  moves the new decision hyperplane in the right direction with respect to  $y_M$

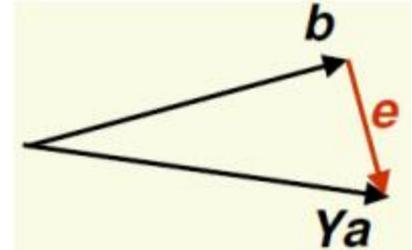
$$(\mathbf{a}^{(k)})^t \mathbf{y}_M \leq 0$$



# MSE Criterion Function

- Minimum squared error approach: find  $\mathbf{a}$  which minimizes the length of the error vector  $\mathbf{e}$

$$\mathbf{e} = \mathbf{Y}\mathbf{a} - \mathbf{b}$$



- Thus minimize the **minimum squared error criterion** function:

$$\mathbf{J}_s(\mathbf{a}) = \|\mathbf{Y}\mathbf{a} - \mathbf{b}\|^2 = \sum_{i=1}^n (\mathbf{a}^t \mathbf{y}_i - b_i)^2$$

- Unlike the perceptron criterion function, we can optimize the minimum squared error criterion function analytically by setting the gradient to 0

# Computing the Gradient

$$\mathbf{J}_s(\mathbf{a}) = \|\mathbf{Y}\mathbf{a} - \mathbf{b}\|^2 = \sum_{i=1}^n (\mathbf{a}^t \mathbf{y}_i - b_i)^2$$

$$\begin{aligned} \nabla \mathbf{J}_s(\mathbf{a}) &= \begin{bmatrix} \frac{\partial \mathbf{J}_s}{\partial \mathbf{a}_0} \\ \vdots \\ \frac{\partial \mathbf{J}_s}{\partial \mathbf{a}_d} \end{bmatrix} = \frac{d\mathbf{J}_s}{d\mathbf{a}} = \sum_{i=1}^n \frac{d}{d\mathbf{a}} (\mathbf{a}^t \mathbf{y}_i - b_i)^2 \\ &= \sum_{i=1}^n 2(\mathbf{a}^t \mathbf{y}_i - b_i) \frac{d}{d\mathbf{a}} (\mathbf{a}^t \mathbf{y}_i - b_i) \\ &= \sum_{i=1}^n 2(\mathbf{a}^t \mathbf{y}_i - b_i) \mathbf{y}_i \\ &= 2\mathbf{Y}^t (\mathbf{Y}\mathbf{a} - \mathbf{b}) \end{aligned}$$

# Pseudo-Inverse Solution

$$\nabla J_s(\mathbf{a}) = 2\mathbf{Y}^t(\mathbf{Y}\mathbf{a} - \mathbf{b})$$

- Setting the gradient to 0:

$$2\mathbf{Y}^t(\mathbf{Y}\mathbf{a} - \mathbf{b}) = 0 \Rightarrow \mathbf{Y}^t\mathbf{Y}\mathbf{a} = \mathbf{Y}^t\mathbf{b}$$

- The matrix  $\mathbf{Y}^t\mathbf{Y}$  is square (it has  $d + 1$  rows and columns) and it is often non-singular
- If  $\mathbf{Y}^t\mathbf{Y}$  is non-singular, its inverse exists and we can solve for  $\mathbf{a}$  uniquely:

$$\mathbf{a} = \boxed{(\mathbf{Y}^t\mathbf{Y})^{-1}\mathbf{Y}^t} \mathbf{b}$$

*pseudo inverse of  $\mathbf{Y}$*   
$$((\mathbf{Y}^t\mathbf{Y})^{-1}\mathbf{Y}^t)\mathbf{Y} = (\mathbf{Y}^t\mathbf{Y})^{-1}(\mathbf{Y}^t\mathbf{Y}) = \mathbf{I}$$

# Ho-Kashyap Procedure

$$\mathbf{J}_{HK}(\mathbf{a}, \mathbf{b}) = \|\mathbf{Y}\mathbf{a} - \mathbf{b}\|^2$$

- As usual, take partial derivatives w.r.t.  $\mathbf{a}$  and  $\mathbf{b}$

$$\nabla_{\mathbf{a}} \mathbf{J}_{HK} = 2\mathbf{Y}^t(\mathbf{Y}\mathbf{a} - \mathbf{b}) = \mathbf{0}$$

$$\nabla_{\mathbf{b}} \mathbf{J}_{HK} = -2(\mathbf{Y}\mathbf{a} - \mathbf{b}) = \mathbf{0}$$

- Use modified gradient descent procedure to find a minimum of  $\mathbf{J}_{HK}(\mathbf{a}, \mathbf{b})$
- Alternate the two steps below until convergence:
  - ① Fix  $\mathbf{b}$  and minimize  $\mathbf{J}_{HK}(\mathbf{a}, \mathbf{b})$  with respect to  $\mathbf{a}$
  - ② Fix  $\mathbf{a}$  and minimize  $\mathbf{J}_{HK}(\mathbf{a}, \mathbf{b})$  with respect to  $\mathbf{b}$

# LDF Summary

- Perceptron procedures
  - Find a separating hyperplane in the linearly separable case,
  - Do not converge in the non-separable case
  - Can force convergence by using a decreasing learning rate, but are not guaranteed a reasonable stopping point
- MSE procedures
  - Converge in separable and not separable case
  - May not find separating hyperplane even if classes are linearly separable
  - Use pseudoinverse if  $\mathbf{Y}^t\mathbf{Y}$  is not singular and not too large
  - Use gradient descent (Widrow–Hoff procedure) otherwise
- Ho–Kashyap procedures
  - always converge
  - find separating hyperplane in the linearly separable case
  - more costly

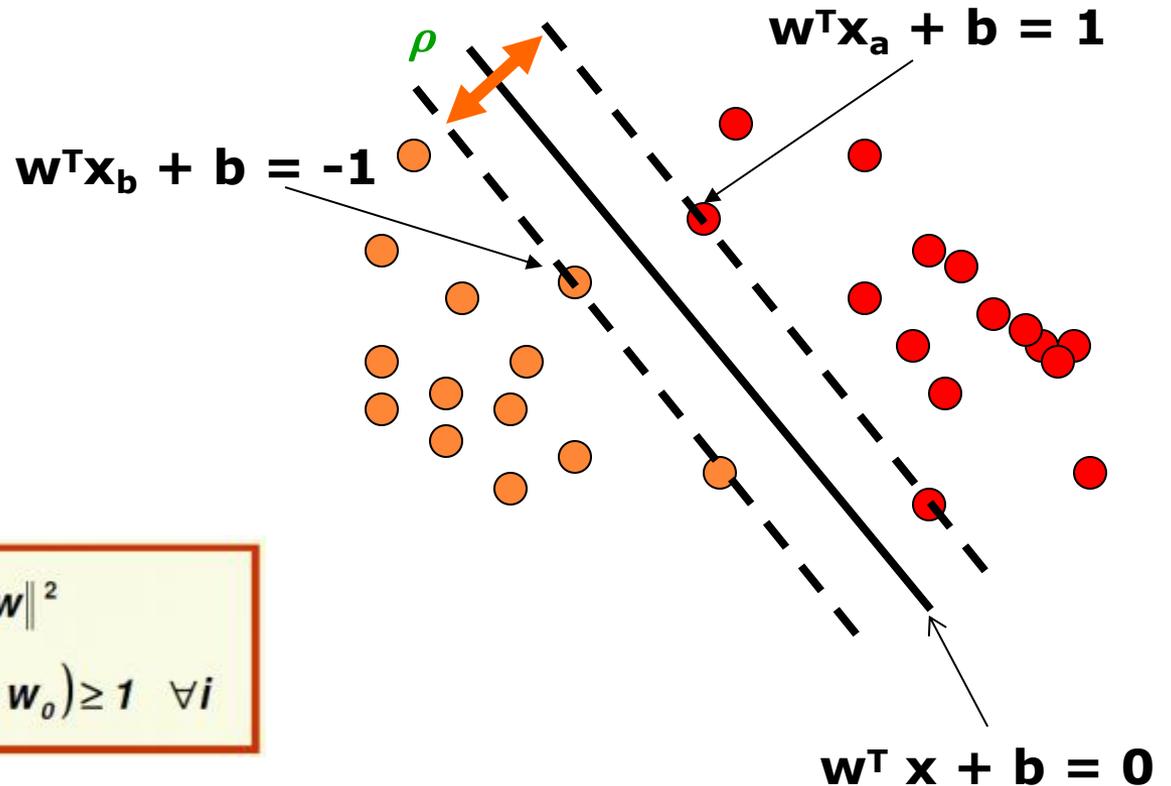
# Linear Support Vector Machine (SVM)

- Maximize margin

$$\rho = \|x_a - x_b\| = 2/\|w\|$$

- Primal problem

$$\begin{aligned} &\text{minimize } J(w) = \frac{1}{2} \|w\|^2 \\ &\text{constrained to } z_i(w^T x_i + w_0) \geq 1 \quad \forall i \end{aligned}$$



# SVM solution: Lagrange multipliers

- We can then swap 'max' and 'min':

$$\begin{aligned} & \min_{\mathbf{w}} \left\{ \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^n \max_{\alpha_j > 0} \alpha_j (1 - z_i(\mathbf{w}^t \mathbf{x}_i + w_0)) \right\} \\ &= \min_{\mathbf{w}} \max_{\alpha_j > 0} \left\{ \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^n \alpha_j (1 - z_i(\mathbf{w}^t \mathbf{x}_i + w_0)) \right\} \\ &= \max_{\alpha_j > 0} \min_{\mathbf{w}} \left\{ \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^n \alpha_j (1 - z_i(\mathbf{w}^t \mathbf{x}_i + w_0)) \right\} \end{aligned}$$

- We can find the optimal  $\mathbf{w}$  as a function of  $\{\alpha_i\}$  by setting the derivatives to zero:

$$\frac{\partial}{\partial \mathbf{w}} J(\mathbf{w}, \alpha) = \mathbf{w} - \sum_{i=1}^n \alpha_i z_i \mathbf{x}_i = \mathbf{0}$$

$$\frac{\partial}{\partial w_0} J(\mathbf{w}, \alpha) = - \sum_{i=1}^n \alpha_i z_i = 0$$

# SVM: Optimal Hyperplane

- Use Kuhn-Tucker Theorem (KTT) condition to convert our problem to:

$$\begin{array}{ll} \text{maximize} & L_D(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j z_i z_j \mathbf{x}_i^t \mathbf{x}_j \\ \text{constrained to} & \alpha_i \geq 0 \quad \forall i \quad \text{and} \quad \sum_{i=1}^n \alpha_i z_i = 0 \end{array}$$

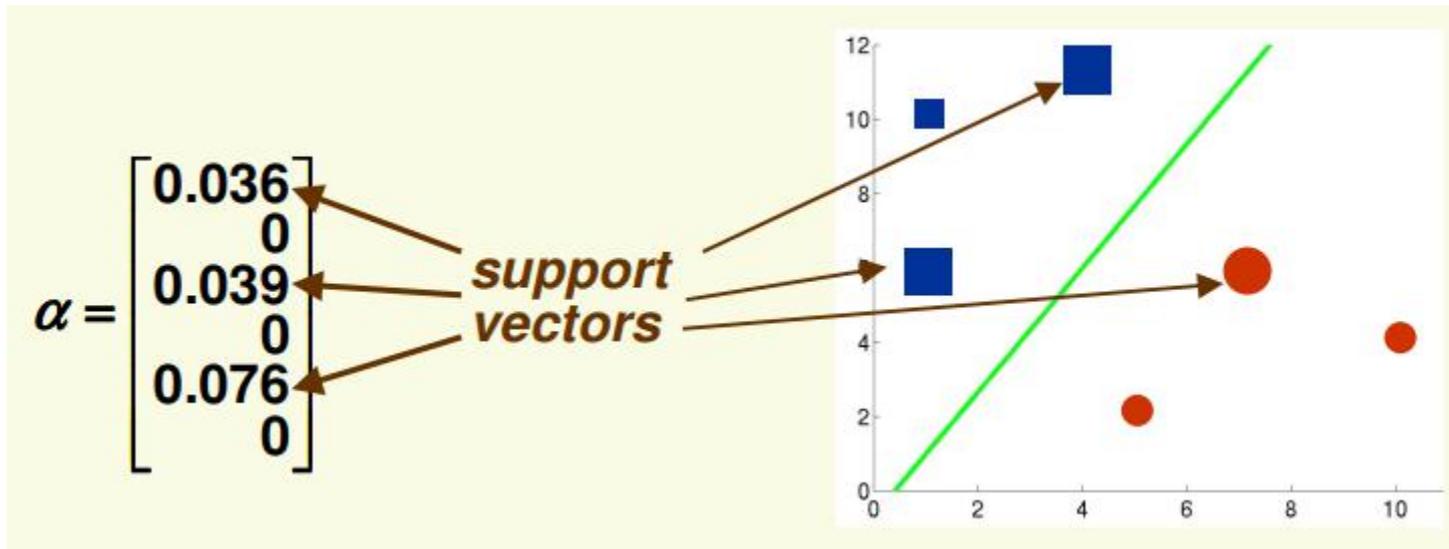
- $\alpha = \{\alpha_1, \dots, \alpha_n\}$  are new variables, one for each sample
- Optimized by quadratic programming

# SVM: Optimal Hyperplane

- After finding the optimal  $a = \{a_1, \dots, a_n\}$
- Final discriminant function:

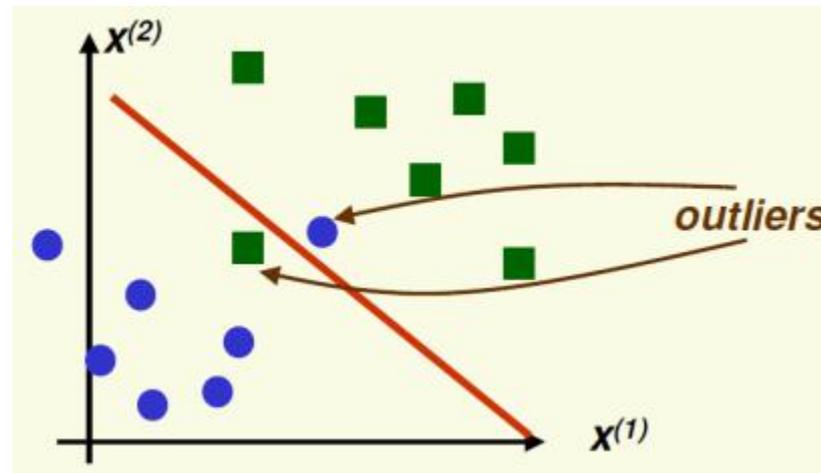
$$g(x) = \left( \sum_{x_i \in S} \alpha_i z_i x_i \right)^t x + w_0 \quad \alpha_i [z_i (w^t x_i + w_0) - 1] = 0 \quad \Rightarrow \quad w_0 = \frac{1}{z_i} - w^t x_i$$

- where  $S$  is the set of **support vectors**



# SVM: Non-Separable Case

- Data are most likely to be not linearly separable, but linear classifier may still be appropriate



- Can apply SVM in non linearly separable case
- Data should be “almost” linearly separable for good performance

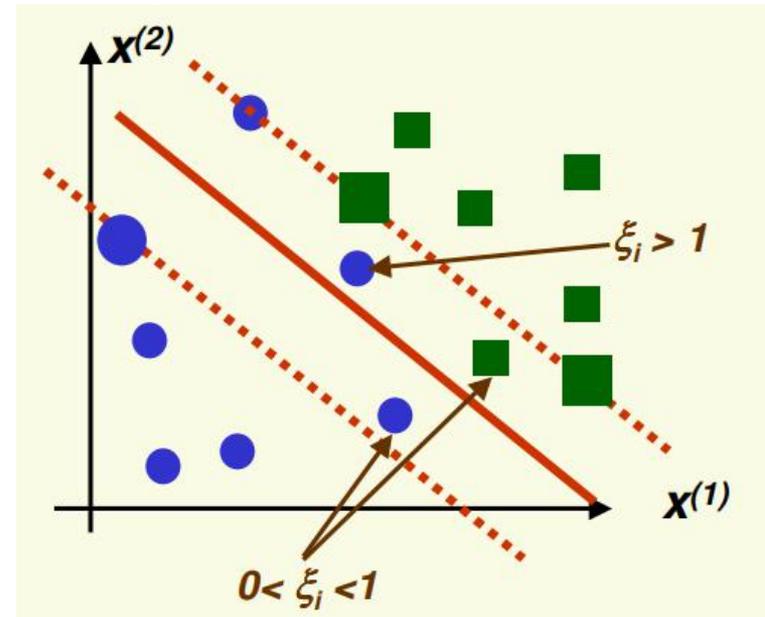
# SVM: Non-Separable Case

- Use slack variables  $\xi_1, \dots, \xi_n$  (one for each sample)

- Change constraints:

$$z_i(w^t x_i + w_0) \geq 1 \quad \forall i \quad \Rightarrow \quad z_i(w^t x_i + w_0) \geq 1 - \xi_i \quad \forall i$$

- $\xi_i$  is a measure of deviation from the ideal for  $x_i$ 
  - $\xi_i \geq 1$ :  $x_i$  is on the wrong side of the separating hyperplane
  - $0 < \xi_i < 1$ :  $x_i$  is on the right side of separating hyperplane but within the region of maximum margin
  - $\xi_i < 0$ : is the ideal case for  $x_i$



# SVM: Non-Separable Case

- Unfortunately this minimization problem is NP-hard due to the discontinuity of  $I(\xi_i)$
- Instead, we minimize

$$J(\mathbf{w}, \xi_1, \dots, \xi_n) = \frac{1}{2} \|\mathbf{w}\|^2 + \beta \sum_{i=1}^n \xi_i$$

*a measure of  
# of misclassified  
examples*

subject to

$$\begin{cases} \mathbf{z}_i (\mathbf{w}^t \mathbf{x}_i + w_0) \geq 1 - \xi_i & \forall i \\ \xi_i \geq 0 & \forall i \end{cases}$$

# SVM: Non-Separable Case

- Use Kuhn–Tucker Theorem (KTT) condition to convert to:

$$\begin{aligned} \text{maximize} \quad & L_D(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j z_i z_j x_i^t x_j \\ \text{constrained to} \quad & 0 \leq \alpha_i \leq \beta \quad \forall i \quad \text{and} \quad \sum_{i=1}^n \alpha_i z_i = 0 \end{aligned}$$

- $w$  and  $w_0$  is computed using:

$$w = \sum_{i=1}^n \alpha_i z_i x_i \quad \Rightarrow \quad \alpha_i [z_i (w^t x_i + w_0) - 1] = 0$$

- Remember that:

$$g(x) = \left( \sum_{x_i \in S} \alpha_i z_i x_i \right)^t x + w_0$$

# What about multi-class SVMs?

- Unfortunately, there is no “definitive” multi-class SVM formulation
- In practice, we have to obtain a multi-class SVM by combining multiple two/ class SVMs
- One vs. others
  - Training: learn an SVM for each class vs. the others
  - Testing: apply each SVM to test example and assign to it the class of the SVM that returns the highest decision value
- One vs. one
  - Training: learn an SVM for each pair of classes
  - Testing: each learned SVM “votes” for a class to assign to the test example

# The Classifiers We Have Learned So Far

Bayesian classifiers {  
MLE classifier  
MAP classifier  
Naive Bayes classifier

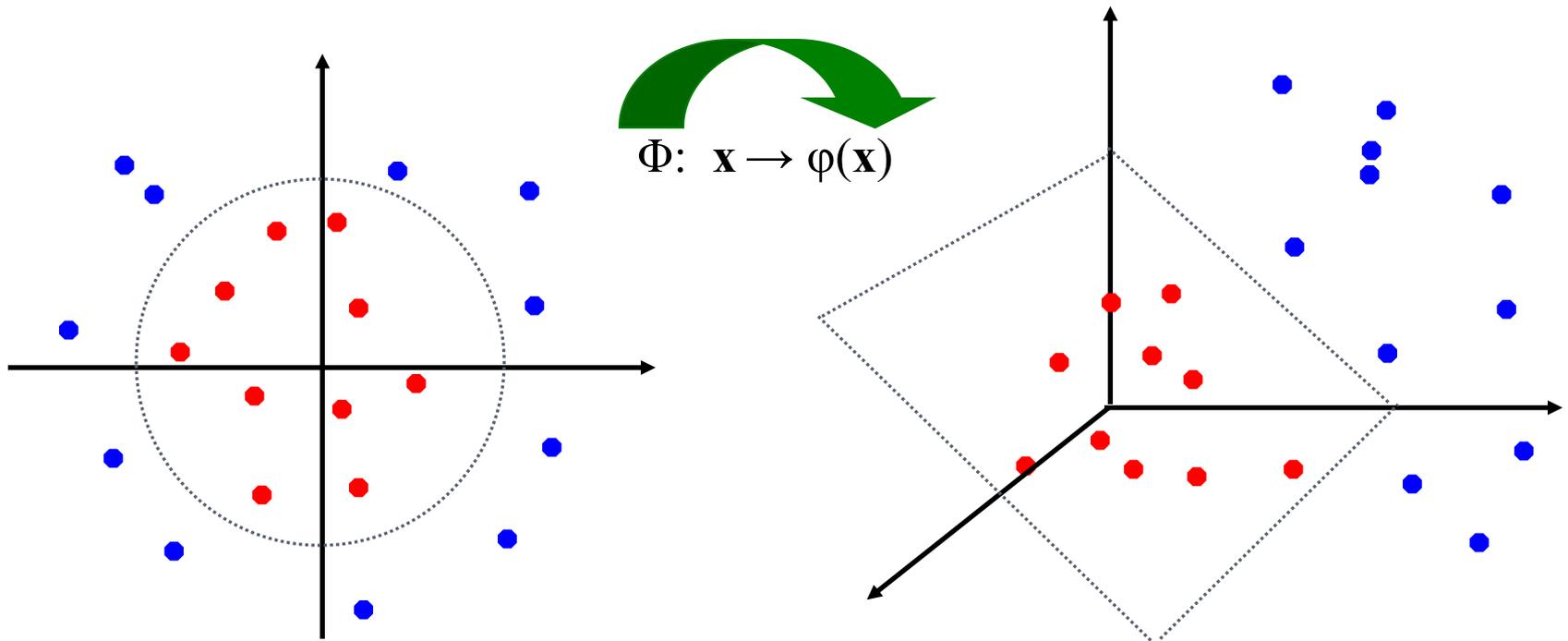
Nonparametric classifiers → KNN classifier

Linear classifiers → LDF (Perceptron rule & Minimu Square Error rule & Ho-Kashyap Procedure) → SVM classifier

Nonlinear classifiers  $\xrightarrow[\text{Feature Mapping } \Phi]{\text{Kernel Tricks}}$  Linear classifiers

# Non-linear SVMs: Feature spaces

- General idea: the original feature space can always be mapped to some higher dimensional feature space where the training set is separable:



# Kernels

- SVM optimization:

$$\text{Maximize } L_D(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j z_i z_j \mathbf{x}_i^t \mathbf{x}_j$$

- Note this optimization depends on samples  $\mathbf{x}_i$  only through the dot product  $\mathbf{x}_i^t \mathbf{x}_j$
- If we lift  $\mathbf{x}_i$  to high dimension using  $\phi(\mathbf{x}_i)$ , we need to compute high dimensional product  $\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$

$$L_D(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j z_i z_j \underbrace{\phi(\mathbf{x}_i)^t \phi(\mathbf{x}_j)}_{K(\mathbf{x}_i, \mathbf{x}_j)}$$

- Idea: find kernel function  $K(\mathbf{x}_i, \mathbf{x}_j)$  s.t.  $K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^t \phi(\mathbf{x}_j)$

# Kernel Trick

- Then we only need to compute  $K(\mathbf{x}_i, \mathbf{x}_j)$  instead of  $\phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j)$
- “**kernel trick**”: do not need to perform operations in high dimensional space explicitly

# Choice of Kernel

- How to choose kernel function  $K(\mathbf{x}_i, \mathbf{x}_j)$ ?
  - $K(\mathbf{x}_i, \mathbf{x}_j)$  should correspond to  $\phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j)$  in a higher dimensional space
  - Mercer's condition tells us which kernel function can be expressed as dot product of two vectors
  - If  $K$  and  $K'$  are kernels  $aK+bK'$  is a kernel
- The mappings  $\phi(\mathbf{x}_i)$  never have to be computed!!

# Pattern recognition design cycle



How can we learn the rule from data?

- **Supervised learning**: a teacher provides a category label or cost for each pattern in the training set.
- **Unsupervised learning**: the system forms clusters or natural groupings of the input patterns.
- **Reinforcement learning**: no desired category is given but the teacher provides feedback to the system such as the decision is right or wrong.

# Learning Algorithms

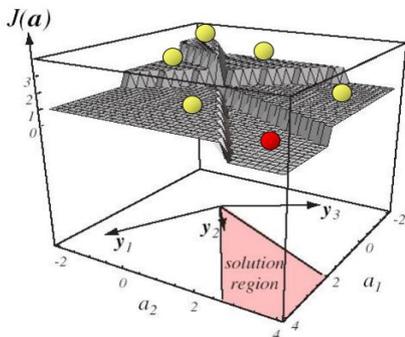
- To design a learning algorithm, we face the following problems:
  - ① Whether to stop ?
  - ② In what direction to proceed ?
  - ③ How long a step to take ?

**Is the criterion satisfactory?**

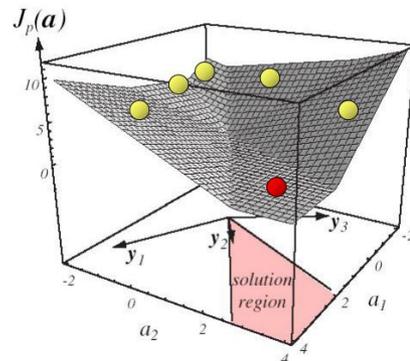
# Criterion Function

- To facilitate learning, we usually define a scalar *critierion function*.
- It usually represents the *penalty* or *cost* of a solution.
- Our goal is to *minimize* its value, i.e., *Function optimization*.

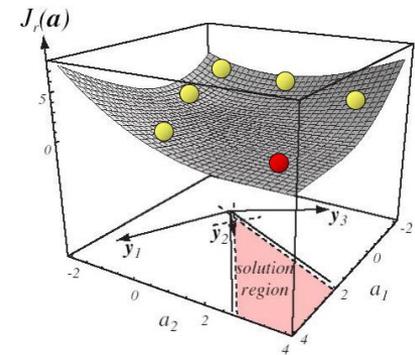
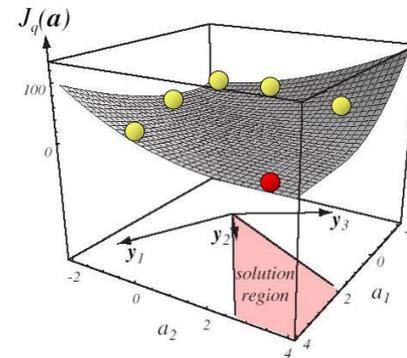
$$J_p(\mathbf{a}) = \sum_{\mathbf{y} \in \mathcal{Y}} (-\mathbf{a}^T \mathbf{y})$$



$$J_q(\mathbf{a}) = \sum_{\mathbf{y} \in \mathcal{Y}} (\mathbf{a}^T \mathbf{y})^2$$



$$J_r(\mathbf{a}) = \sum_{\mathbf{y} \in \mathcal{Y}} \frac{(\mathbf{a}^T \mathbf{y} - b)^2}{\|\mathbf{y}\|^2}$$

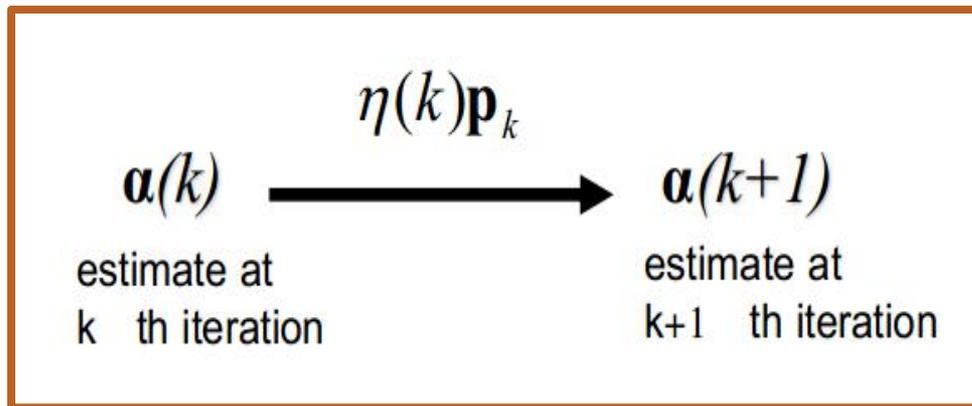


# Learning Using Iterative Optimization

- Minimize an error function  $J(\alpha)$  (e.g., classification error) with respect to  $\alpha$ :
- **Minimize**  $J(\alpha)$  iteratively:  $\alpha(k+1) = \alpha(k) + \eta(k)\mathbf{p}_k$

learning rate

search direction



How should we choose  $\mathbf{p}_k$ ?

# Choosing $\mathbf{p}_k$ using Gradient Descent

$$\mathbf{p}_k = -\nabla J(\mathbf{a}(k))$$

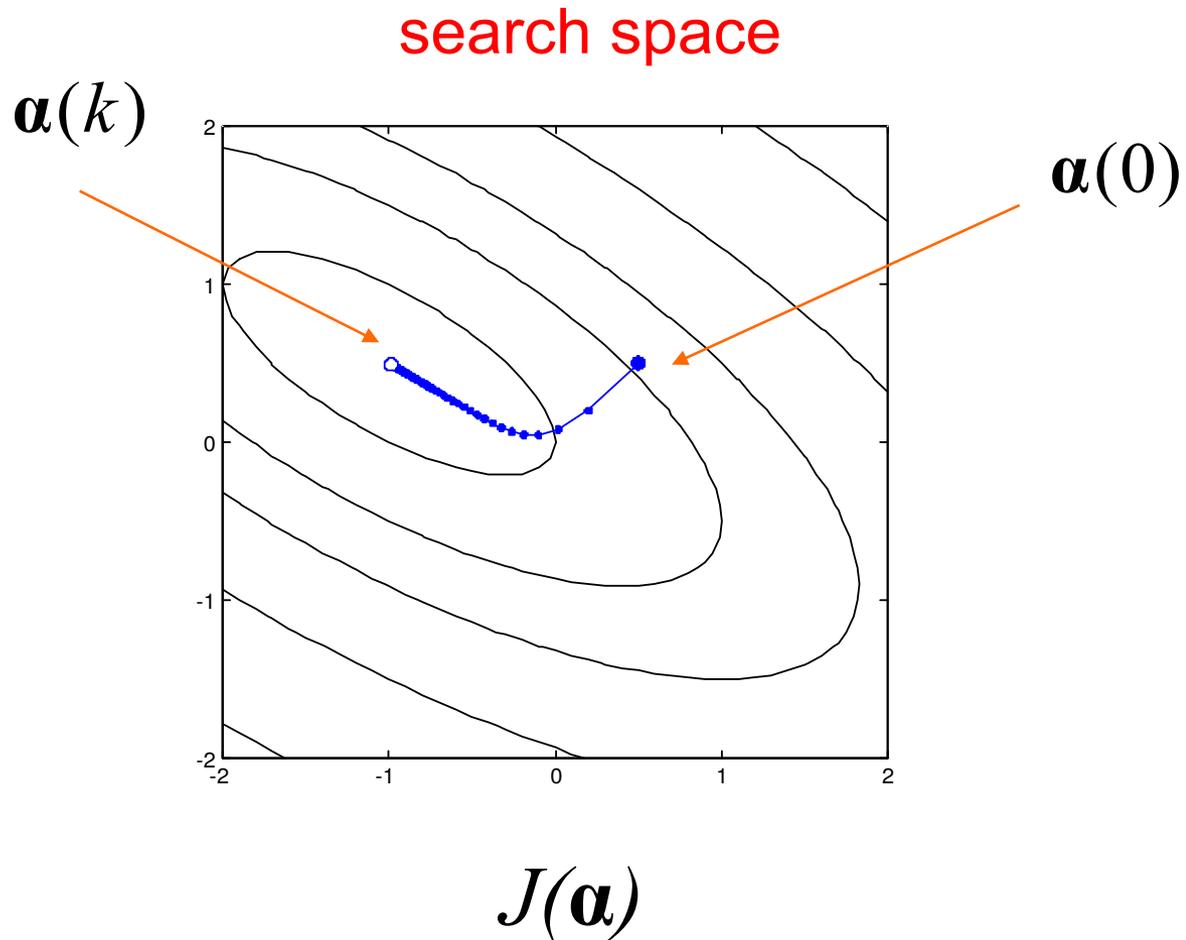
$$\mathbf{a}(k + 1) = \mathbf{a}(k) - \eta(k) \nabla J(\mathbf{a}(k))$$

## Algorithm 1 (Basic gradient descent)

```
1 begin initialize  $\mathbf{a}$ , criterion  $\theta$ ,  $\eta(\cdot)$ ,  $k = 0$   
2   do  $k \leftarrow k + 1$   
3      $\mathbf{a} \leftarrow \mathbf{a} - \eta(k) \nabla J(\mathbf{a})$   
4   until  $\eta(k) \nabla J(\mathbf{a}) < \theta$   
5 return  $\mathbf{a}$   
6 end
```

( $\mathbf{a} = \alpha$ )

# Gradient Decent (cont'd)



# Gradient Decent (cont'd)

- How to choose the learning rate  $\eta(k)$ ?

Taylor series approximation

( $\mathbf{a} = \boldsymbol{\alpha}$ )

$$J(\mathbf{a}) \simeq J(\mathbf{a}(k)) + \nabla J^t (\mathbf{a} - \mathbf{a}(k)) + \frac{1}{2} (\mathbf{a} - \mathbf{a}(k))^t \mathbf{H} (\mathbf{a} - \mathbf{a}(k))$$

Hessian (2<sup>nd</sup> derivatives)

Setting  $\mathbf{a} = \mathbf{a}(k+1)$  and using  $\mathbf{a}(k+1) = \mathbf{a}(k) - \eta(k) \nabla J(\mathbf{a}(k))$

$$J(\mathbf{a}(k+1)) \simeq J(\mathbf{a}(k)) - \eta(k) \|\nabla J\|^2 + \frac{1}{2} \eta^2(k) \nabla J^t \mathbf{H} \nabla J$$

$$\eta(k) = \frac{\|\nabla J\|^2}{\nabla J^t \mathbf{H} \nabla J} \quad \text{optimum learning rate}$$

# Choosing $\mathbf{p}_k$ using Newton's Method

$$\mathbf{p}_k = -\mathbf{H}^{-1}\nabla J(\mathbf{a}(k))$$

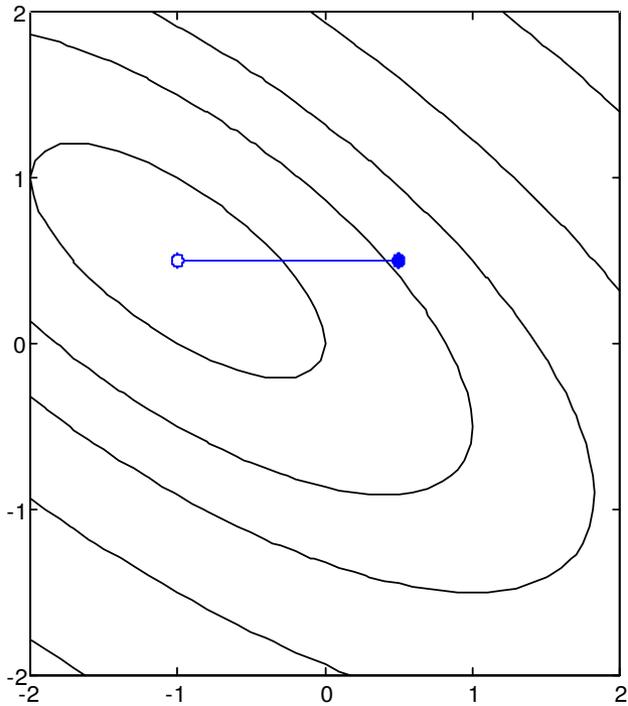
$$\mathbf{a}(k+1) = \mathbf{a}(k) - \mathbf{H}^{-1}\nabla J \text{ requires inverting } \mathbf{H}$$

## Algorithm 2 (Newton descent)

```
1 begin initialize a, criterion  $\theta$ 
2   do
3      $\mathbf{a} \leftarrow \mathbf{a} - \mathbf{H}^{-1}\nabla J(\mathbf{a})$    ( $\mathbf{a} = \boldsymbol{\alpha}$ )
4     until  $\mathbf{H}^{-1}\nabla J(\mathbf{a}) < \theta$ 
5   return a
6 end
```

# Newton's Method (cont'd)

$J(\alpha)$

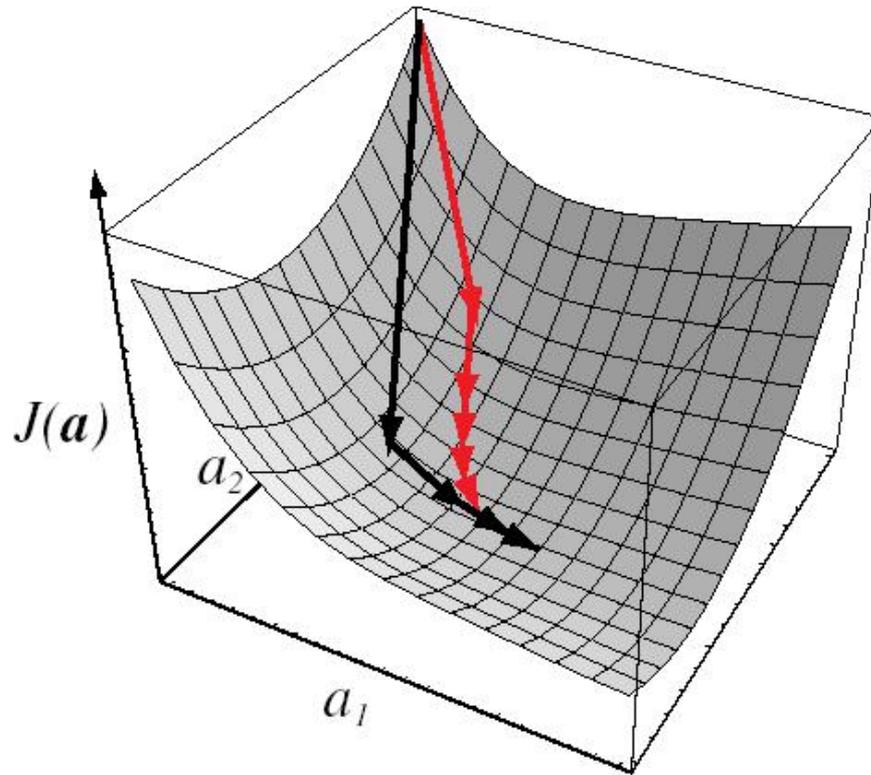


If  $J(\alpha)$  is **quadratic**,  
Newton's method  
converges in **one step!**

# Gradient decent vs Newton's method

Newton's method

Gradient Decent



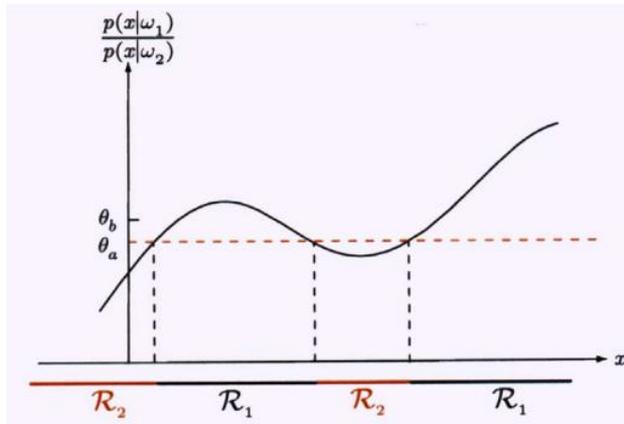
# Pattern recognition design cycle



- How can we estimate the performance with training samples?
- How can we predict the performance with future data?
- Problems of overfitting and generalization.

# Receiver Operating Characteristic (ROC) Curve

- Every classifier typically employs some kind of a threshold.

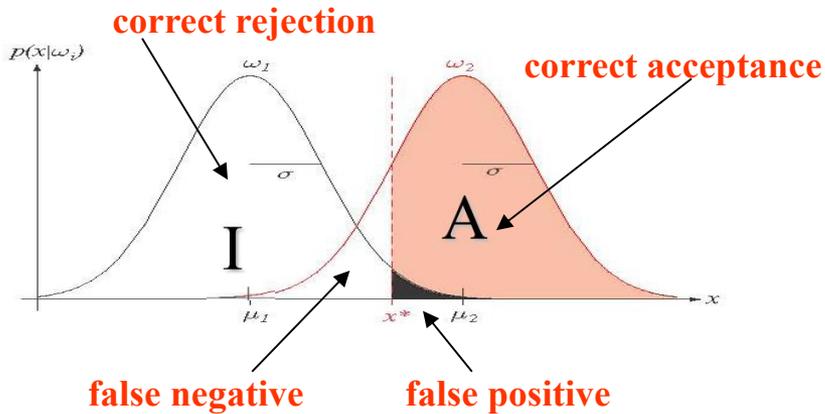


$$\theta_a = P(\omega_2) / P(\omega_1)$$

$$\theta_b = \frac{P(\omega_2)(\lambda_{12} - \lambda_{22})}{P(\omega_1)(\lambda_{21} - \lambda_{11})}$$

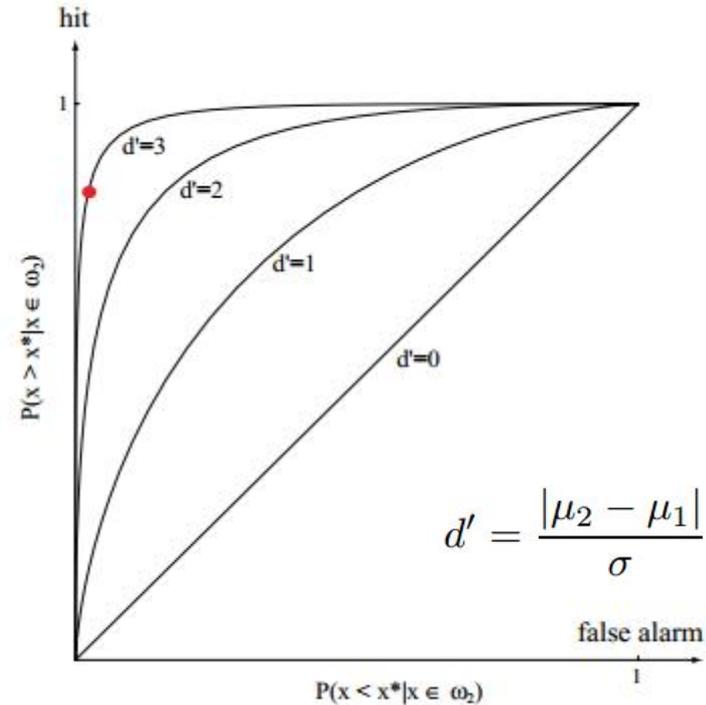
- Changing the threshold will affect the performance of the classifier.
- ROC curves allow us to evaluate the performance of a classifier using **different** thresholds.

# ROC Curve



$$TPR = TP/P = TP/(TP + FN)$$

$$FPR = FP/N = FP/(FP + TN)$$



**FPR:** False Positive Rate (X-axis)  
**TRR:** True Positive Rate (Y-axis)

# Overfitting

- Prediction error: probability of test pattern not in class with max posterior (**true**)
- Training error: probability of test pattern not in class with max posterior (**estimated**)
- Classifier optimized w.r.t. training error
  - Training error: optimistically biased estimate of prediction error

# Overfitting

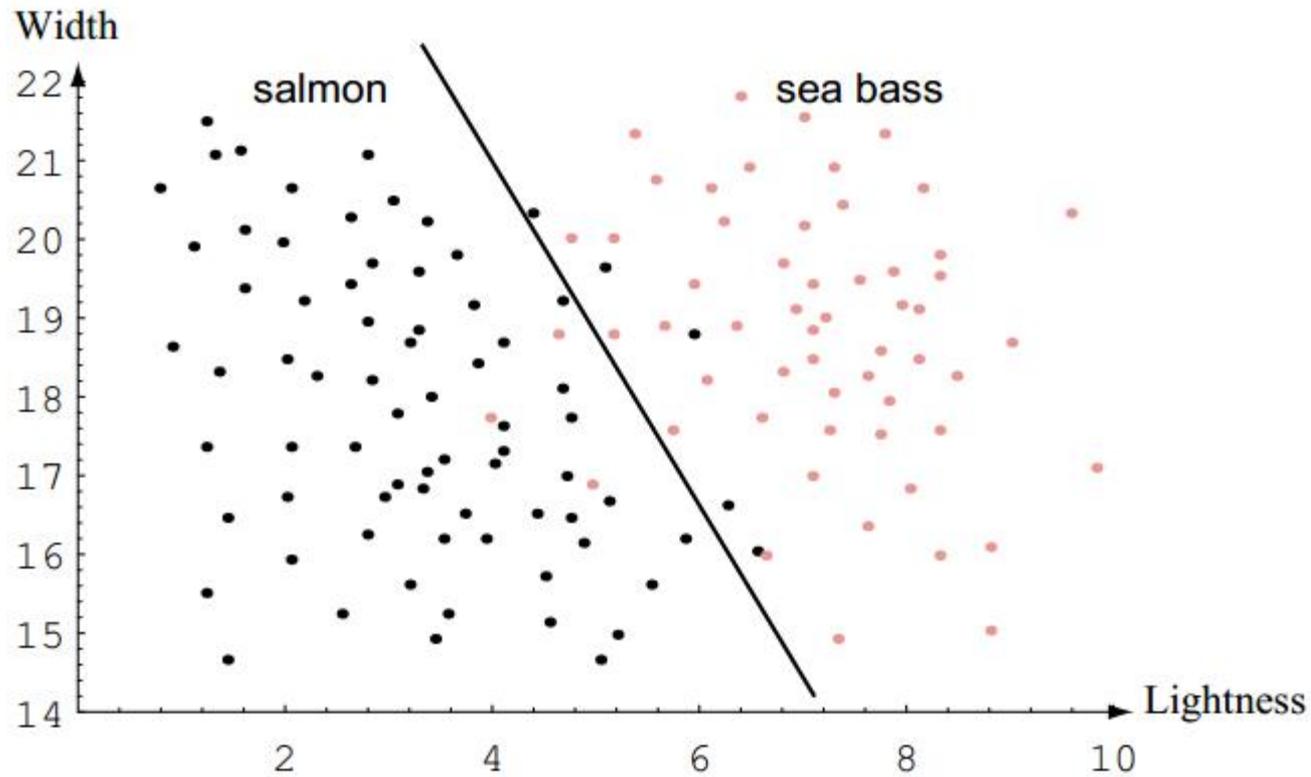
- Overfitting: a learning algorithm overfits the training data if it outputs a solution  $w$  when another solution  $w'$  exists such that:

$$\text{error}_{\text{train}}(w) < \text{error}_{\text{train}}(w')$$

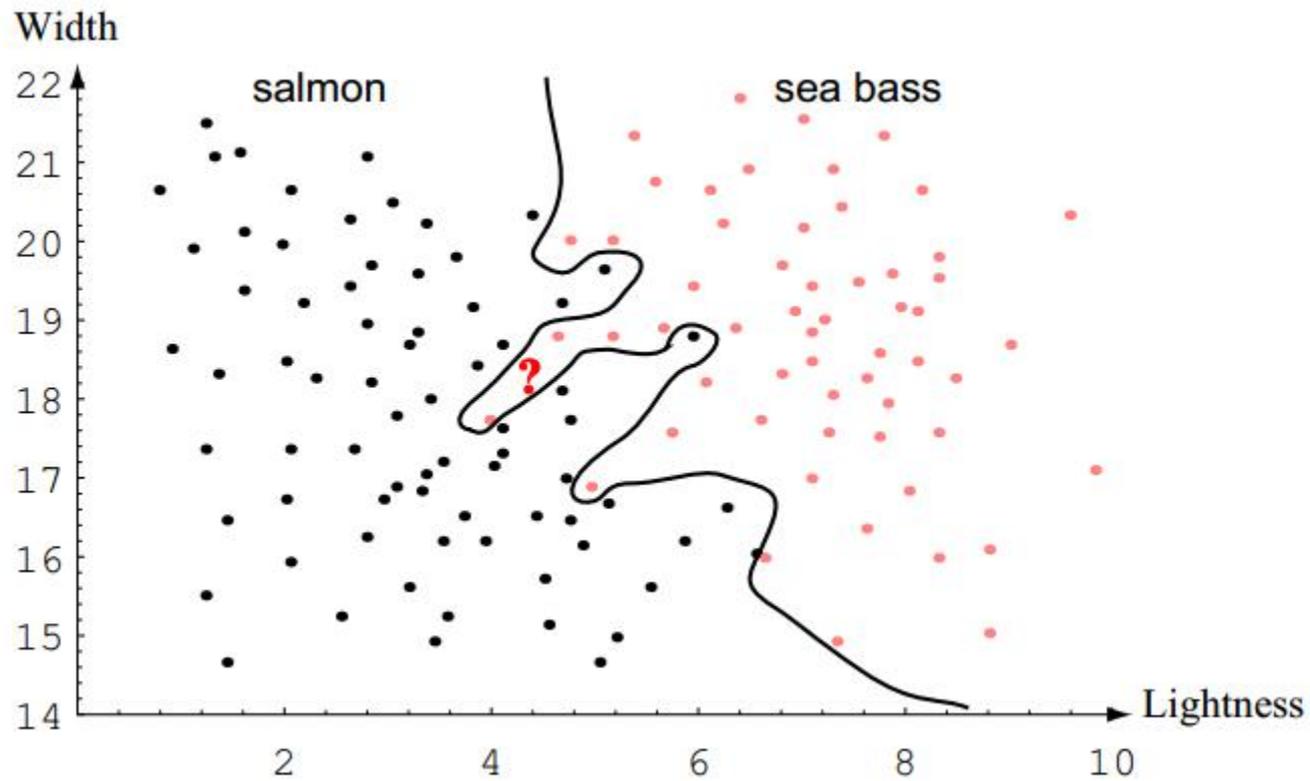
AND

$$\text{error}_{\text{true}}(w') < \text{error}_{\text{true}}(w)$$

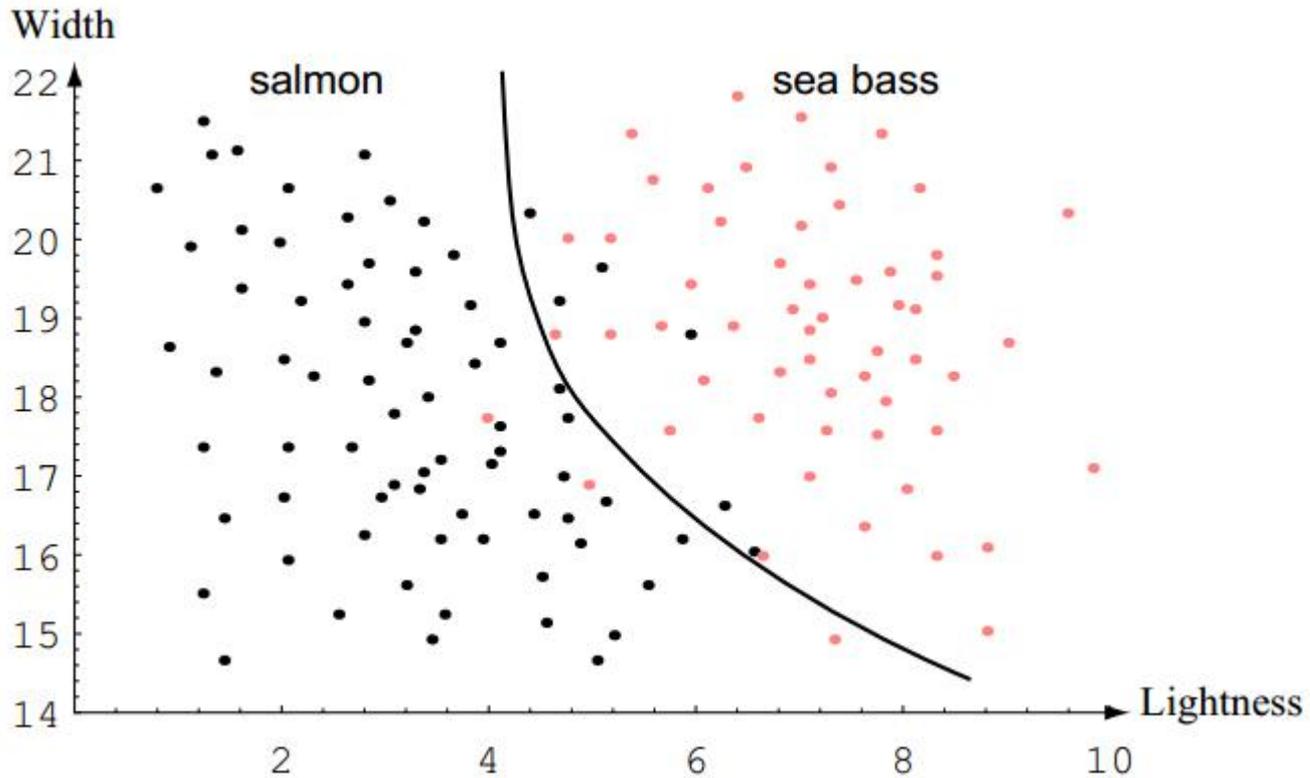
# Example: Fish Classifier



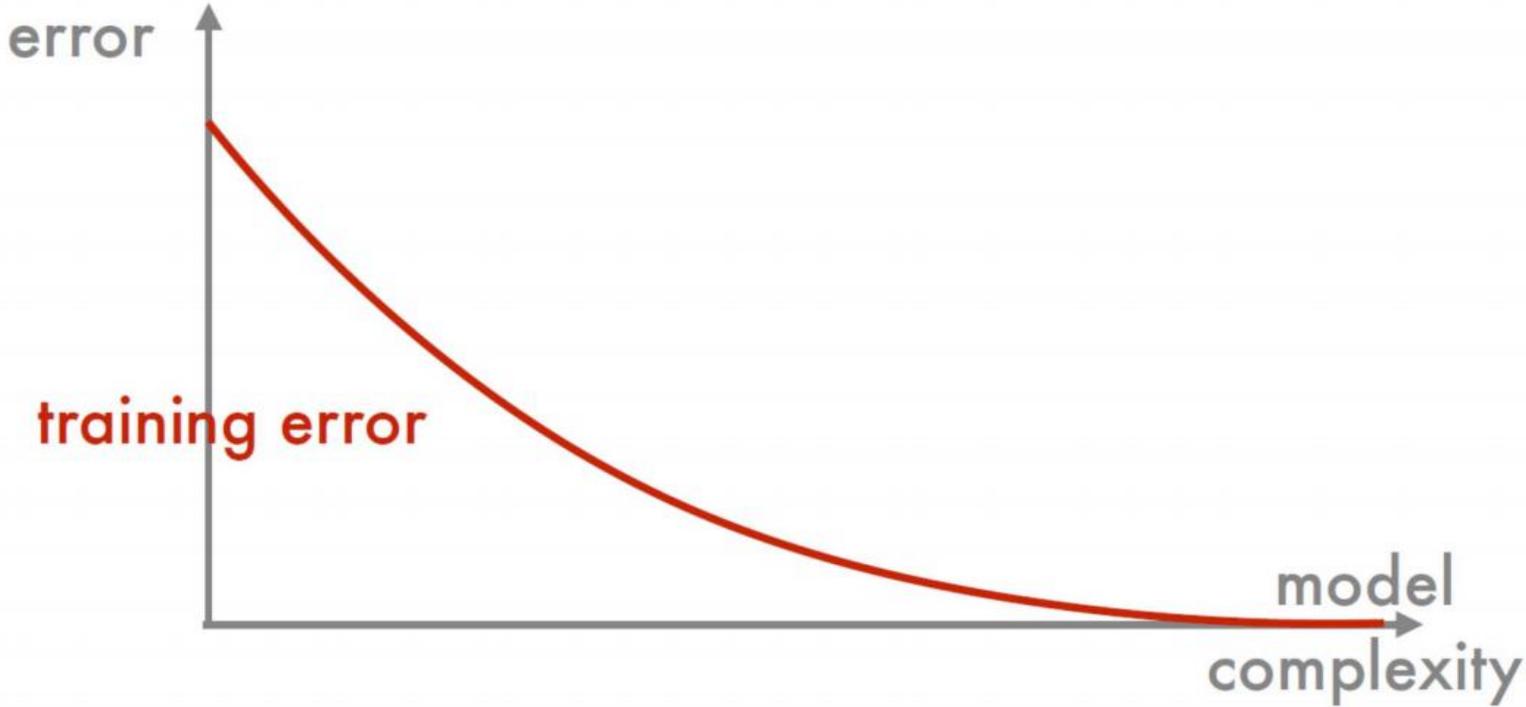
# Minimum Training Error



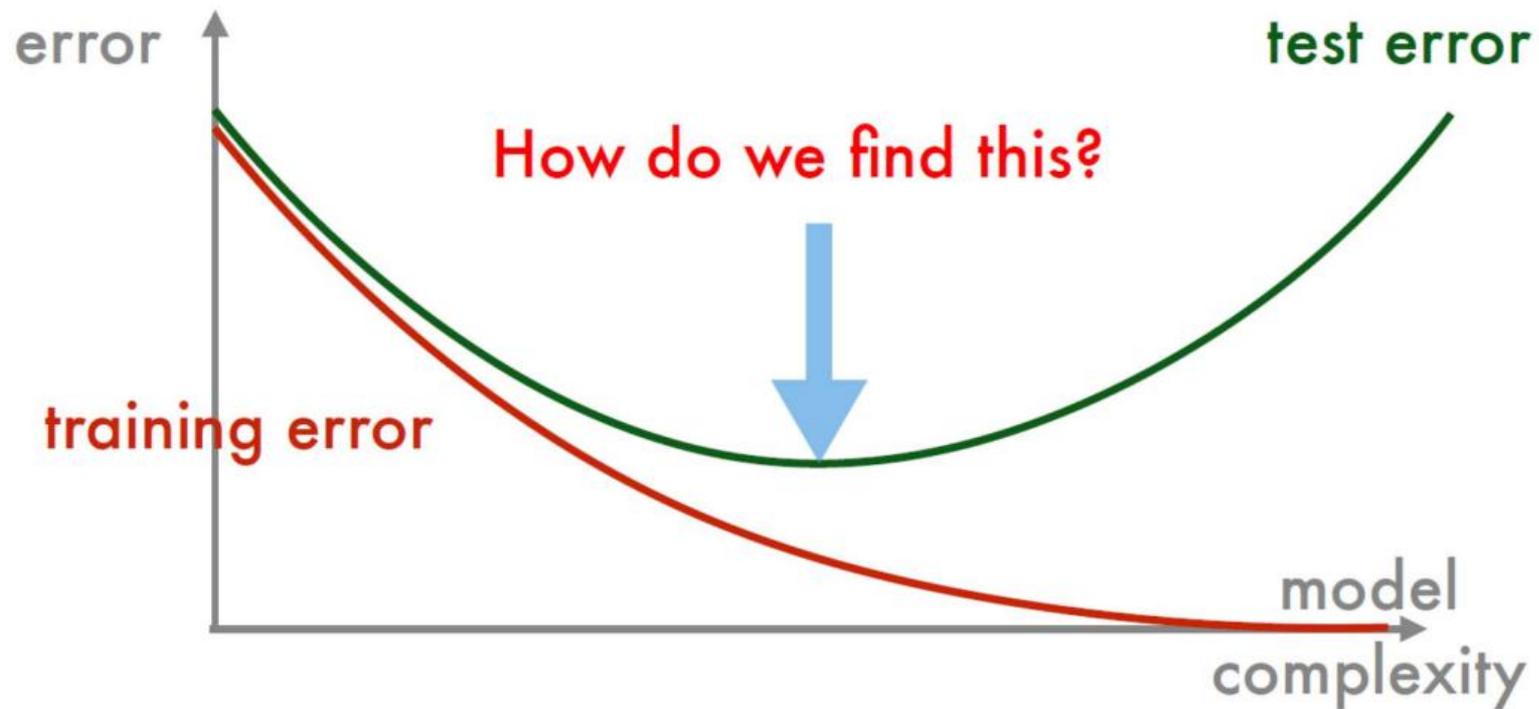
# Final Decision Boundary



# Typical Behavior

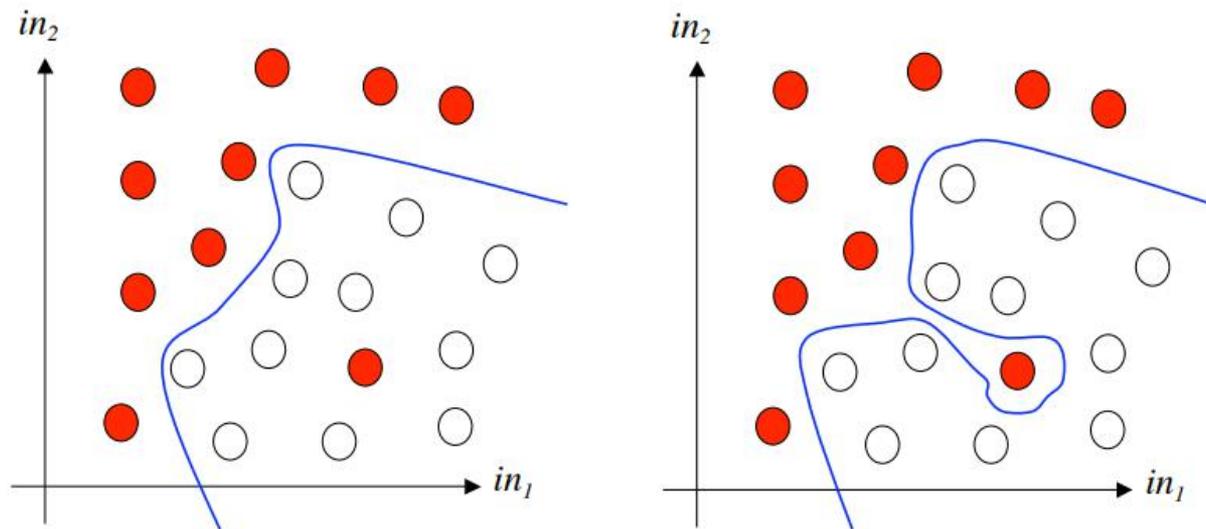


# Typical Behavior



# Typical Behavior

- The aim is to get a classification model to generalize to classify new inputs appropriately.
- If the training data is known to contain noise, we don't necessarily want the training data to be classified totally accurately, because that is likely to reduce the generalization ability.



# *Q & A*