

Classification: Linear Discriminant Functions

CE-725: Statistical Pattern Recognition
Sharif University of Technology
Spring 2013

Soleymani

Outline

- ▶ Discriminant functions
 - ▶ Linear Discriminant functions + BW: nonlinear discriminant functions
- ▶ Linear Discriminant Function
 - ▶ Least Mean Squared Error Method LMS alg.
 - ▶ Sum of Squared Error Method SSE
 - ▶ Perceptron Error Correction method
- ▶ Multi-class problems
 - ▶ Linear machine
 - ▶ Completely Linearly Separable
 - ▶ Pairwise Linearly Separable
- ▶ Generalized LDFs

Classification Problem

- ▶ **Given: Training set**
 - ▶ labeled set of N input-output pairs $D = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$
 - ▶ $y \in \{1, \dots, c\}$
- ▶ **Goal: Given an input \mathbf{x} , assign it to one of c classes**

Types of Classifiers

- ▶ Probabilistic classification approaches (previous lectures):
 - ▶ Generative
 - ▶ Discriminative
- ▶ Discriminant function
 - ▶ Various procedures for determining discriminant functions (some of them are statistical)
 - ▶ However, they don't require knowledge of the forms of underlying probability distributions

Discriminant Functions

- ▶ **Discriminant functions:** A popular way of representing a classifier

- ▶ A discriminant function $g_i(\mathbf{x})$ for each class ω_i ($i = 1, \dots, c$):
 - ▶ \mathbf{x} is assigned to class ω_i if:

/omega

$$g_i(\mathbf{x}) > g_j(\mathbf{x}) \quad \forall j \neq i$$

- ▶ **Decision surfaces** (boundaries) can also be found using discriminant functions
 - ▶ Boundary of the \mathcal{R}_i and \mathcal{R}_j : $\forall \mathbf{x}, g_i(\mathbf{x}) = g_j(\mathbf{x})$

Probabilistic Discriminant Functions

Prior probability: $p(w_i)$

- ▶ Maximum likelihood

- ▶ $g_i(\mathbf{x}) = p(\mathbf{x}|\omega_i)$ a.k.a. conditional probability: $p(\mathbf{x}|w_i)$,
e.g. If you're positive(w_i), and you tested (\mathbf{x}) positive.

- ▶ Bayesian Classifier

- ▶ $g_i(\mathbf{x}) = P(\omega_i|\mathbf{x})$ a.k.a. Posterior probability: $p(w_i|\mathbf{x})$,
e.g. If you tested (\mathbf{x}) positive, and you're ACTUALLY positive(w_i).

- ▶ Expected Loss (Conditional Risk)

- ▶ $g_i(\mathbf{x}) = -R(a_i|\mathbf{x})$ +BW: a_i : treat as actions

Discriminant Functions: Two-Category

- ▶ For **two-category** problem, we can only find a function $g : \mathbb{R}^d \rightarrow \mathbb{R}$
 - ▶ $g_1(\mathbf{x}) = g(\mathbf{x})$
 - ▶ $g_2(\mathbf{x}) = -g(\mathbf{x})$
- ▶ Decision surface: $g(\mathbf{x}) = 0$
- ▶ First, we explain **two-category** classification problem and then discuss the **multi-category** problems.

Linear Discriminant Functions

▶ Linear Discriminant Functions (LDFs)

assumption ▶ Decision boundaries are linear in \mathbf{x} , or linear in some given set of functions of \mathbf{x}

▶ Why LDFs?

▶ They can be **optimal** for some problems

▶ **E.g.**, if the underlying distributions $p(\mathbf{x}|\omega_j)$ are gaussians having equal covariance

▶ Even when they are not optimal, we can use their simplicity

▶ LDFs are relatively easy to compute

▶ In the absence of information suggesting otherwise, linear classifiers are an attractive candidates for initial, trial classifiers.

LDFs: Two Category

- ▶ $g(\mathbf{x}; \mathbf{w}) = \mathbf{w}^T \mathbf{x} = w_0 + w_1 x_1 + \dots w_d x_d$ vs BW: high order poly forms
 - ▶ $\mathbf{x} = [1 \ x_1 \ x_2 \ \dots \ x_d]$ \longrightarrow Augmented, which is to add one dim, i.e. +1
 - ▶ $\mathbf{w} = [w_0 \ w_1 \ w_2 \ \dots \ w_d]$
 - ▶ w_0 : bias
 - ▶ \mathbf{w} contains the parameters we need to set

if $g(\mathbf{x}; \mathbf{w}) = \mathbf{w}^T \mathbf{x} \geq 0$ then ω_1
else ω_2

- ▶ Decision surface (boundary)^H: $g(\mathbf{x}; \mathbf{w}) = 0$
 - ▶ The equation $g(\mathbf{x}; \mathbf{w}) = 0$ defines the decision surface separating samples of the two categories
 - ▶ When $g(\mathbf{x}; \mathbf{w})$ is linear, the decision surface is a hyperplane. ???

LDFs : Two Category

- ▶ Decision boundary is a $(d - 1)$ -dimensional hyperplane H in the d -dimensional feature space
 - ▶ The orientation of H is determined by the normal vector $[w_1, \dots, w_d]$
 - ▶ w_0 is the location of the surface is determined by the bias.

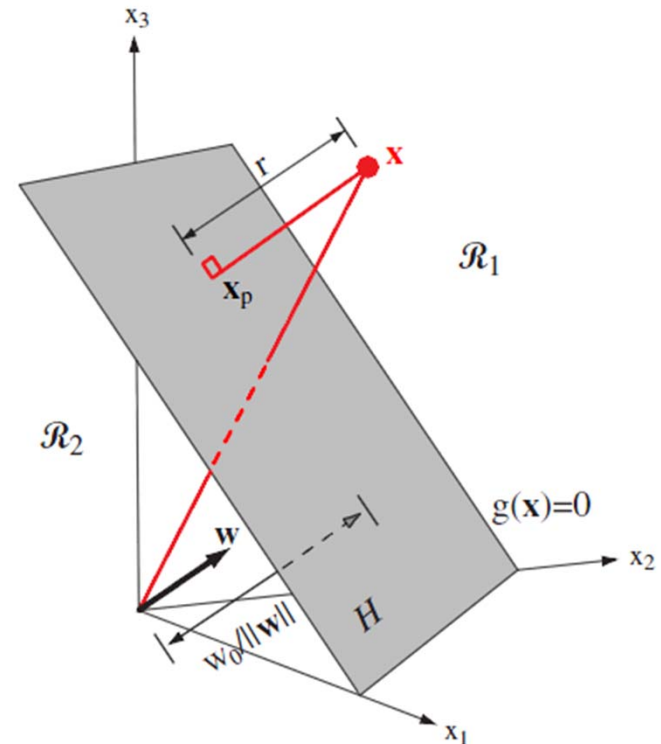
+BW:笔记里有推导

$$\mathbf{x} = \mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|}$$

$$g(\mathbf{x}) = r\|\mathbf{w}\| \Rightarrow r = \frac{g(\mathbf{x})}{\|\mathbf{w}\|}$$

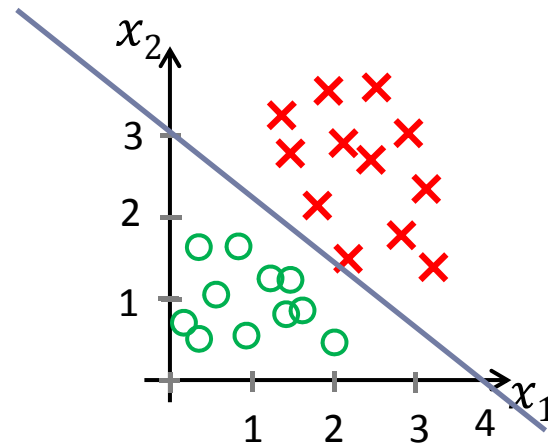


$g(\mathbf{x})$ is proportional to the signed distance from \mathbf{x} to H



LDFs: Two Category

$$3 + \frac{3}{4}x_1 + x_2 = 0$$



if $\mathbf{w}^T \mathbf{x} \geq 0$ ω_1
else ω_2

$$\mathbf{w} = [3, 0.75, 1]$$

$$\mathbf{x} = [1, x_1, x_2]$$

LDFs: Cost Function

- ▶ Finding LDFs is formulated as an optimization problem
 - ▶ A cost function is needed and a procedure is used to solve it.
- ▶ Criterion or cost functions for classification:
 - ▶ Average training error or loss incurred in classifying training samples
 - ▶ A small training error does not guarantee a small test error
 - ▶ We will investigate several cost functions for the classification problem

LDFs: Methods

- ▶ Many classification methods are based on LDFs:
 - ▶ Mean Squared Error i.e. MSE
 - ▶ Sum of Squared Error i.e. SSE
 - ▶ Perceptron i.e. error correction method
 - ▶ Fisher Linear Discriminant Analysis (LDA) [next lectures]
 - ▶ SVM [next lecture]

Main Steps in Methods based on LDFs

- ▶ We have specified the class of discriminant functions as linear
- ▶ Select how to measure the prediction loss
 - ▶ Based on the training set $D = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n$, a cost function $J(\mathbf{w})$ is defined (e.g., SSE) that is a function of the classifier parameters
- ▶ Solve the resulting optimization problem to find parameters:
 - ▶ Find optimal $\hat{g}(\mathbf{x}) = g(\mathbf{x}; \hat{\mathbf{w}})$ where $\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} J(\mathbf{w})$

Bayes vs. LDFs' Cost Function

- ▶ **Bayes minimum error** classifier:

$$\min_{\alpha(\cdot)} E_{\mathbf{x}, y} [L_{0-1}(\alpha(\mathbf{x}), y)]$$

- ▶ If we know the probabilities in advance then the above optimization problem will be solved easily.

- ▶ $\alpha(\mathbf{x}) = \operatorname{argmax}_y P(y|\mathbf{x})$

- ▶ We only have a set of training samples \mathcal{D} instead of $p(\mathbf{x}, y)$ and usually optimize the following problem:

$$\min_{\alpha(\cdot)} \frac{1}{N} \sum_{i=1}^N L(\alpha(\mathbf{x}^{(i)}), y^{(i)})$$

Mean Squared Error

Two-category: $y \in \{-1, 1\}$
 $y = -1$ for ω_2 , $y = 1$ for ω_1

Squared Loss $\text{Loss}(\alpha(\mathbf{x}), y) = (y - \alpha(\mathbf{x}))^2$

$$J(\mathbf{w}) = E_{xy} \left[(y - g(\mathbf{x}; \mathbf{w}))^2 \right] = E[(y - \mathbf{w}^T \mathbf{x})^2] \quad E?$$

Gradient Descent

$$\nabla_{\mathbf{w}} J(\mathbf{w}) = \mathbf{0} \Rightarrow 2E[\mathbf{x}(y - \mathbf{w}^T \mathbf{x})] = \mathbf{0}$$

$$\Rightarrow E[\mathbf{x}y] = E[\mathbf{x}\mathbf{x}^T]\mathbf{w} \Rightarrow \hat{\mathbf{w}} = \frac{E[\mathbf{x}y]}{E[\mathbf{x}\mathbf{x}^T]} = \mathbf{R}_x^{-1} \mathbf{R}_{xy}$$

$$\mathbf{R}_x = \begin{bmatrix} E[x_1 x_1] & \cdots & E[x_1 x_d] \\ \vdots & \ddots & \vdots \\ E[x_d x_1] & \cdots & E[x_d x_d] \end{bmatrix} \quad \mathbf{R}_{xy} = \begin{bmatrix} E[x_1 y] \\ \vdots \\ E[x_d y] \end{bmatrix}$$

Sum of Squared Error (SSE)

- ▶ Cost function: Prediction errors on the training set:
 - ▶ empirical loss on N training samples

$$\begin{aligned} J(\mathbf{w}) &= \frac{1}{N} \sum_{i=1}^N \text{Loss}(y^{(i)}, g(\mathbf{x}^{(i)}; \mathbf{w})) \\ &= \frac{1}{N} \sum_{i=1}^N (y^{(i)} - g(\mathbf{x}^{(i)}; \mathbf{w}))^2 \\ &= \frac{1}{N} \sum_{i=1}^N (y^{(i)} - \mathbf{w}^T \mathbf{x}^{(i)})^2 \end{aligned}$$

Sum of Squared Error (SSE)

$$J(\mathbf{w}) = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2$$

unknown (i.e. \mathbf{w}) is marked at the latter place

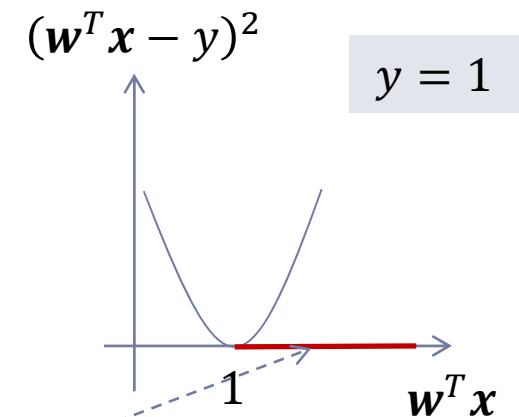
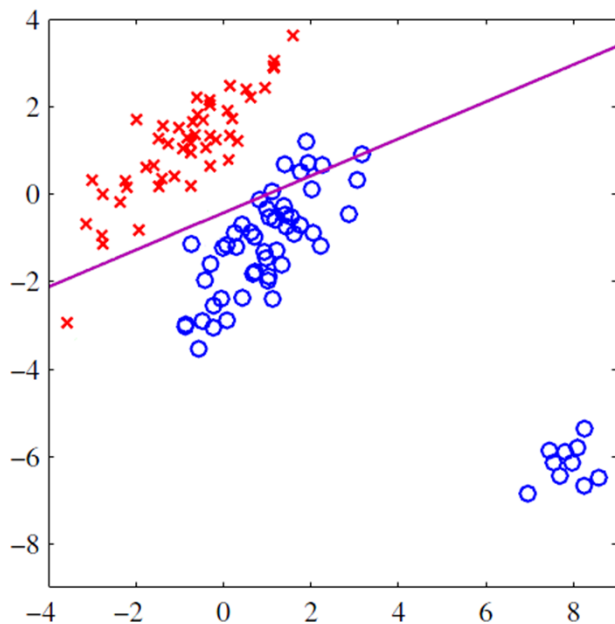
$$\mathbf{X} = \begin{bmatrix} 1 & x_1^{(1)} & \dots & x_d^{(1)} \\ 1 & x_1^{(2)} & \dots & x_d^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(N)} & \dots & x_d^{(N)} \end{bmatrix} \quad \mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(N)} \end{bmatrix}$$

$N \times (d+1)$ column vectors

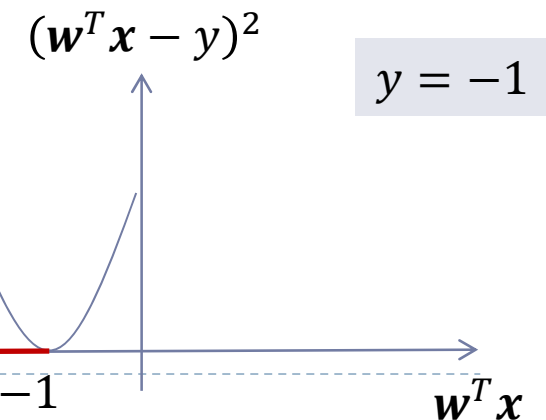
$$\begin{aligned} \nabla_{\mathbf{w}} J(\mathbf{w}) = \mathbf{0} &\Rightarrow 2\mathbf{X}^T(\mathbf{X}\mathbf{w} - \mathbf{y}) = \mathbf{0} \Rightarrow \mathbf{X}^T\mathbf{X}\mathbf{w} = \mathbf{X}^T\mathbf{y} \Rightarrow \hat{\mathbf{w}} \\ &= \underbrace{(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T}_{\text{Pseudo Inverse of } \mathbf{X}} \mathbf{y} \end{aligned}$$

Sum of Squared Error (SSE)

- ▶ SSE penalizes 'too correct' predictions
 - ▶ 'too correct' predictions: samples lie a long way on the correct side of the decision
- ▶ It also lack robustness to noise

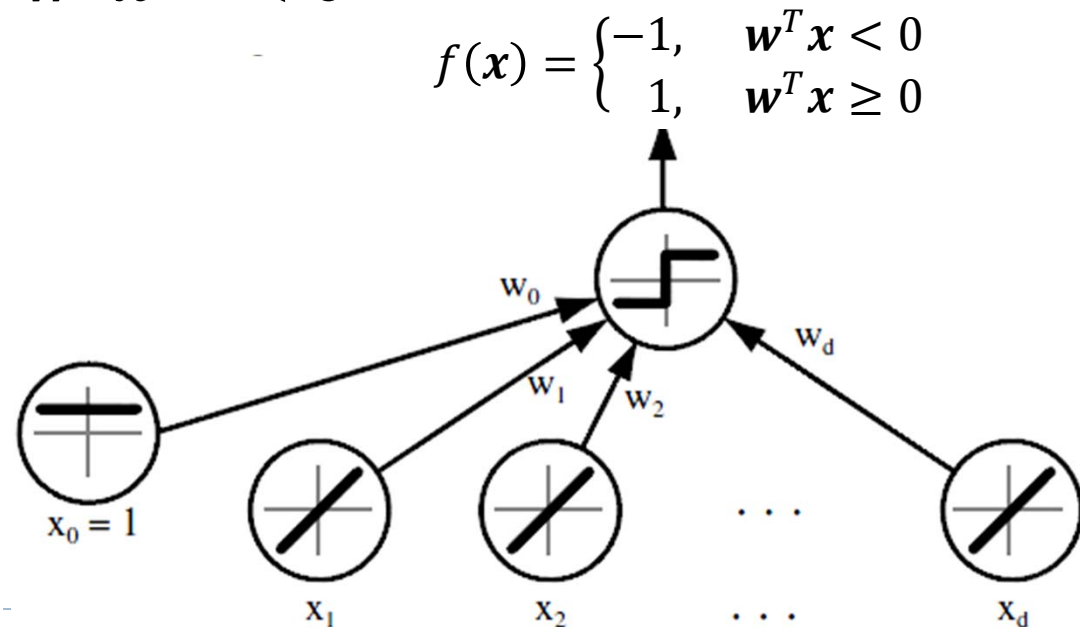


Correct predictions that
are penalized by SSE



Perceptron a.k.a. Error Correction method

- ▶ Two-category: $y \in \{-1, 1\}$
 - ▶ $y = -1$ for ω_2 , $y = 1$ for ω_1
/omega
- ▶ Goal: $\forall i, \mathbf{x}^{(i)} \in \omega_1 \Rightarrow \mathbf{w}^T \mathbf{x}^{(i)} > 0$
 $\forall i, \mathbf{x}^{(i)} \in \omega_2 \Rightarrow \mathbf{w}^T \mathbf{x}^{(i)} < 0$



Perceptron Criterion

- ▶ Only misclassified training samples affect the discriminant functions:

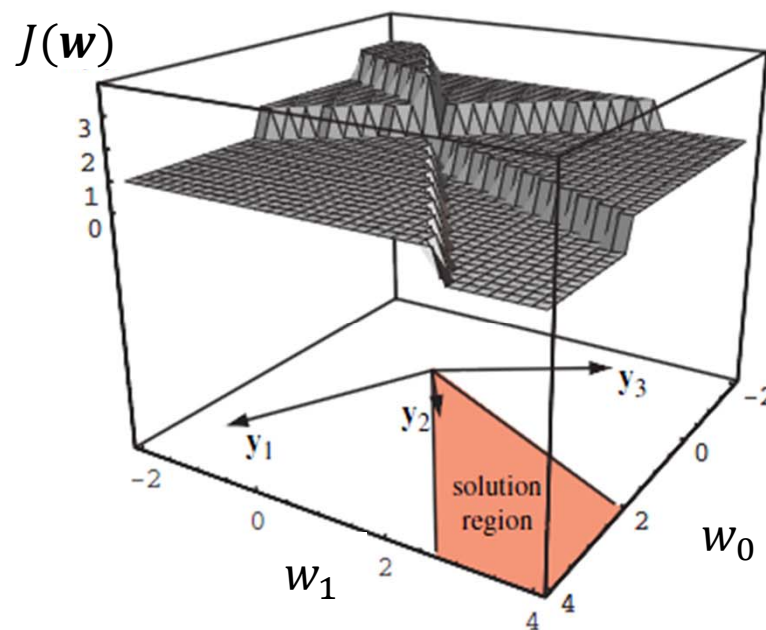
$$J_P(\mathbf{w}) = - \sum_{i \in \mathcal{M}} \mathbf{w}^T \mathbf{x}^{(i)} y^{(i)}$$

\mathcal{M} : subset of training data that are misclassified

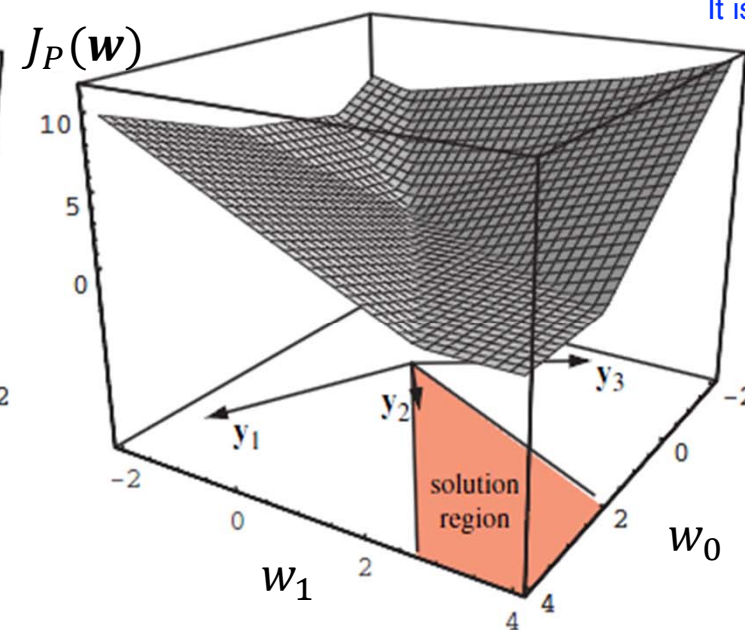
- ▶ Many solutions? Which solution among them?

Perceptron vs. Other Criteria

Here, we are designing a better criteria or cost function.



Misclassification

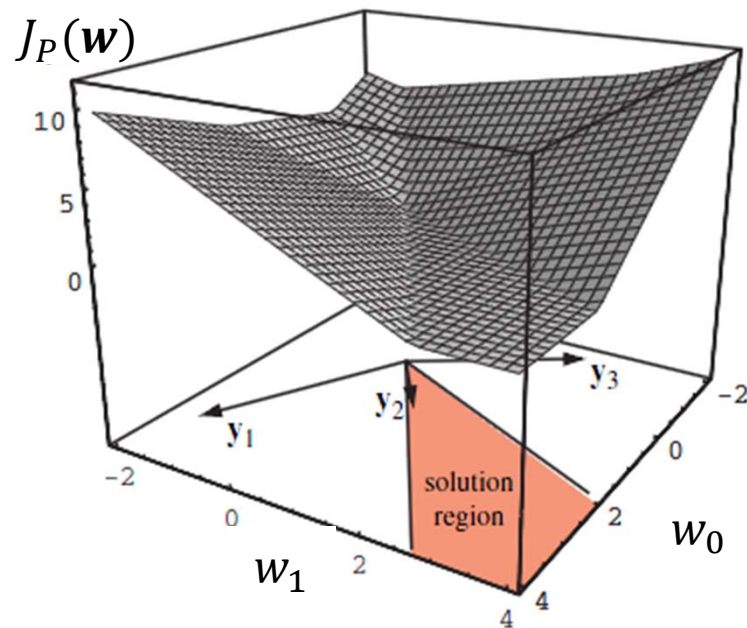


Perceptron

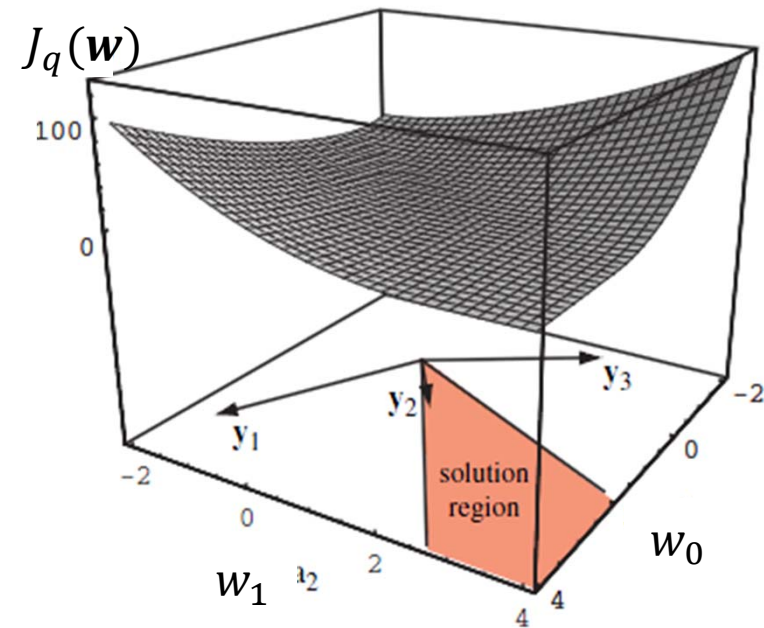
Perceptron $J_P(\mathbf{w})$ is piecewise linear.
It is not continuous.

[Duda,Hart&Stork]

Some classification criteria



Perceptron



SSE

[Duda,Hart&Stork]

Batch Perceptron

“Gradient Descent” to solve the optimization problem:

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \underset{\text{/eta}}{\eta} \underset{\text{/nabla}}{\nabla}_{\mathbf{w}} J_P(\mathbf{w}^t)$$

$$\nabla_{\mathbf{w}} J_P(\mathbf{w}) = - \sum_{i \in \mathcal{M}} \mathbf{x}^{(i)} y^{(i)}$$

Batch Perceptron converges in finite number of steps
for linearly separable data

Initialize $\mathbf{w}, t \leftarrow 0$

Repeat

$$\mathbf{w} = \mathbf{w} + \eta \sum_{i \in \mathcal{M}} \mathbf{x}^{(i)} y^{(i)}$$

$$t \leftarrow t + 1$$

Until $\eta \sum_{i \in \mathcal{M}} \mathbf{x}^{(i)} y^{(i)} < \theta$

Single-Sample Perceptron

- ▶ If $\mathbf{x}^{(i)}$ is misclassified:

$$\mathbf{w}^{t+1} = \mathbf{w}^t + \eta \mathbf{x}^{(i)} y^{(i)}$$

Fixed-Increment Single Sample Perceptron

Initialize $\mathbf{w}, k \leftarrow 0$

repeat

$k \leftarrow (k + 1) \bmod N$

if $\mathbf{x}^{(i)}$ is misclassified then

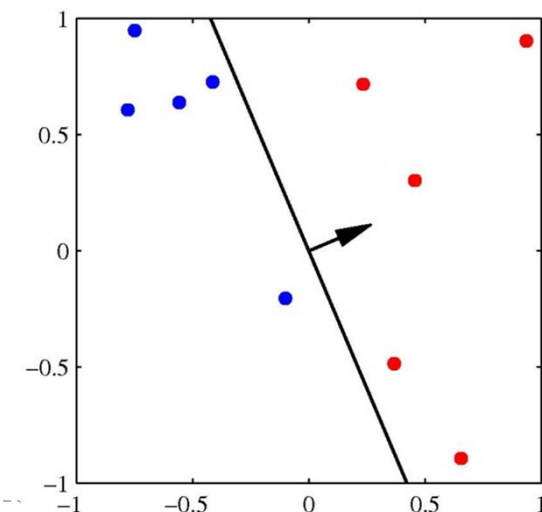
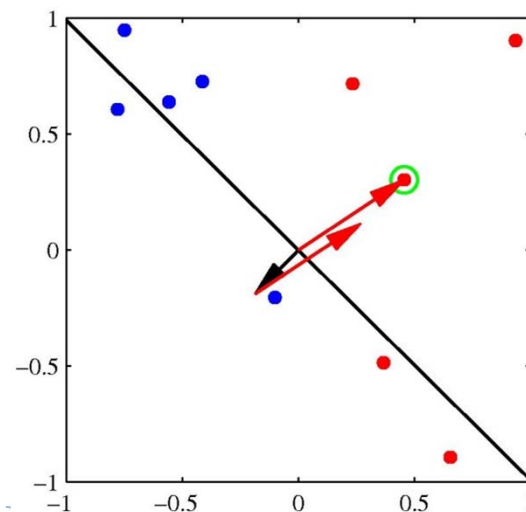
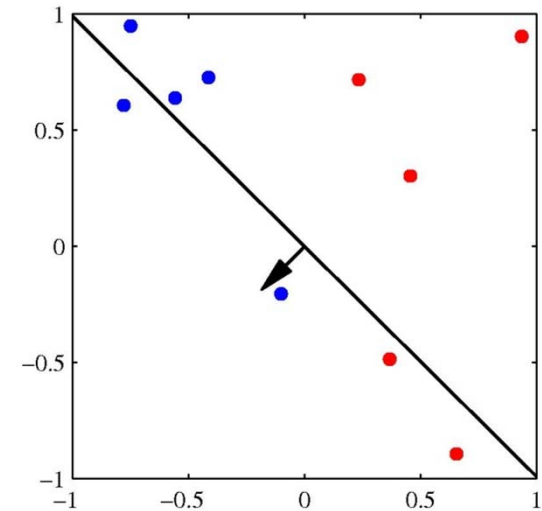
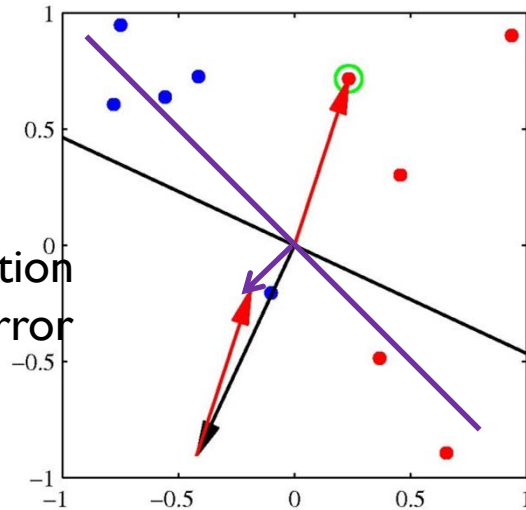
$\mathbf{w} = \mathbf{w} + \mathbf{x}^{(i)} y^{(i)}$

Until all patterns properly classified

η can be set to 1 ←

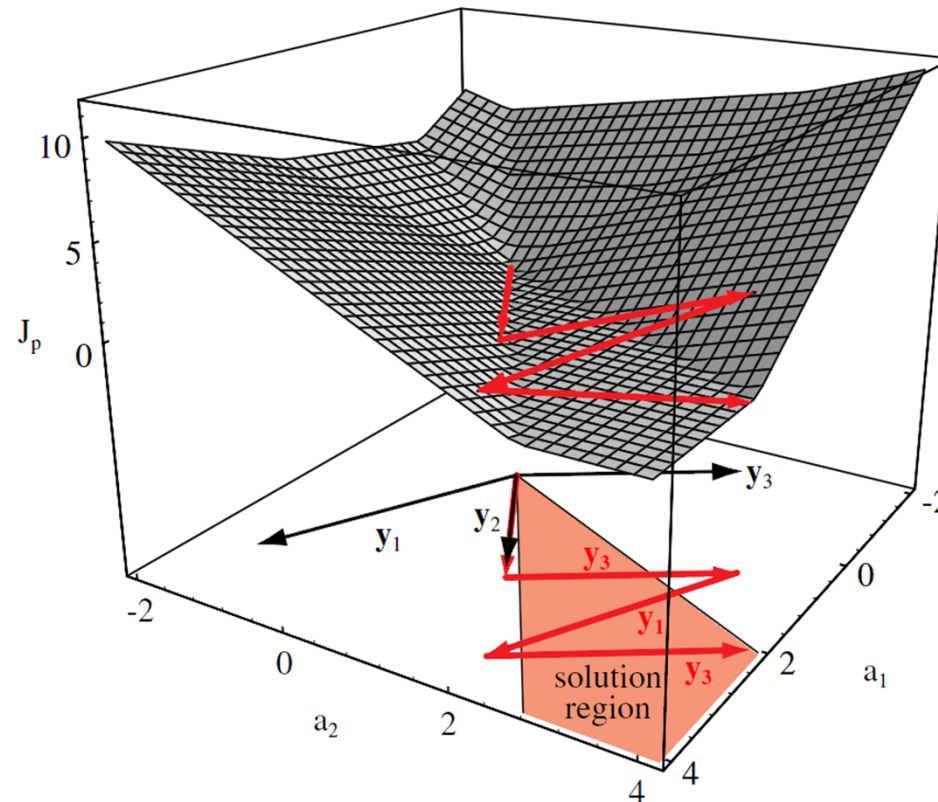
Convergence of Perceptron

Change w in a direction
that corrects the error

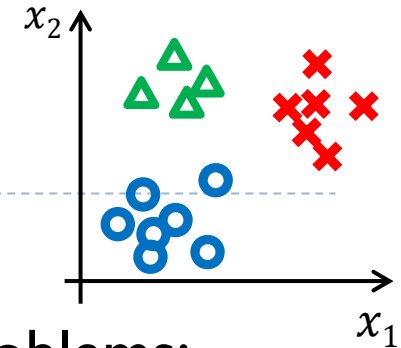


Convergence of Perceptron

- ▶ If the training data set is linearly separable, the single-sample perceptron algorithm is also guaranteed to find a solution in a finite number of steps



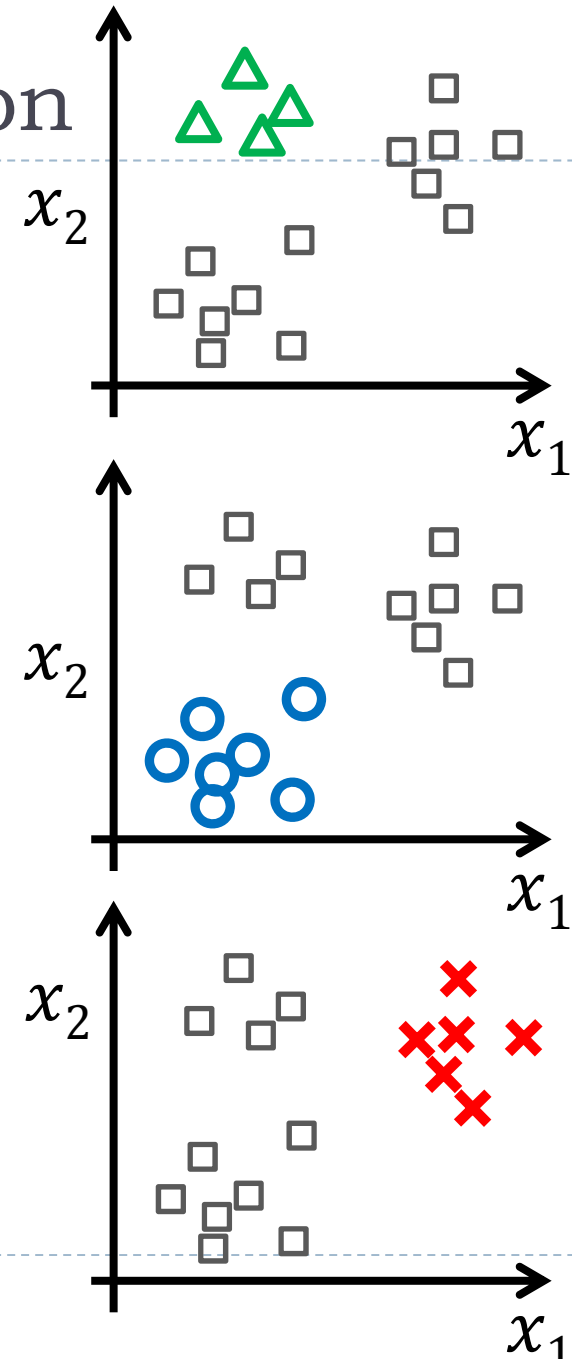
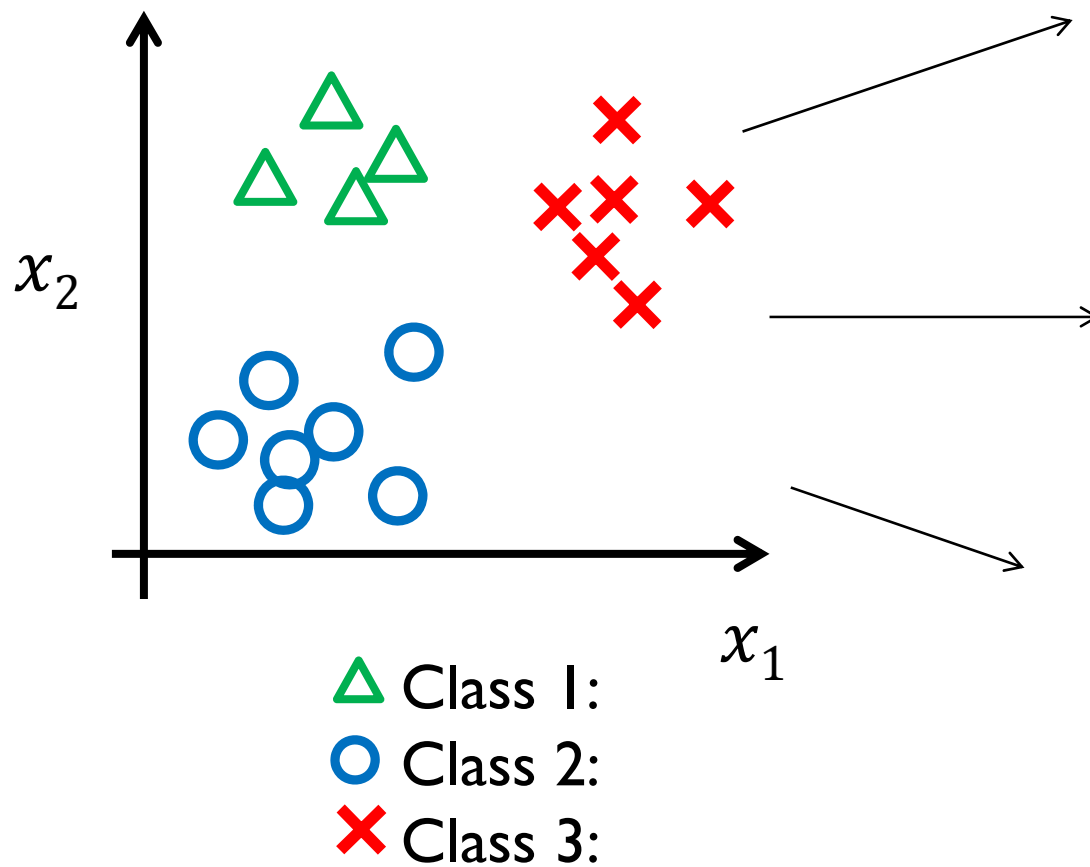
LDFs: Multi-Category



- ▶ Linear discriminant functions for multi-category problems:
 - ▶ Linear machine :
 - ▶ A discriminant function $g_i(x)$ for each class i
 - ▶ Converting the problem to a set of two-class problems:
 - ▶ “one versus rest” or “one against all”
 - For each class ω_i , an LDF separates samples of ω_i from all the other samples.
 - Totally linearly separable
 - ▶ “one versus one”
 - $c(c - 1)/2$ LDFs are used, one to separate samples of a pair of classes.
 - Pairwise linearly separable
- ▶ Converting the problem to a set of two-class problems can lead to **regions in which the classification is undefined.**

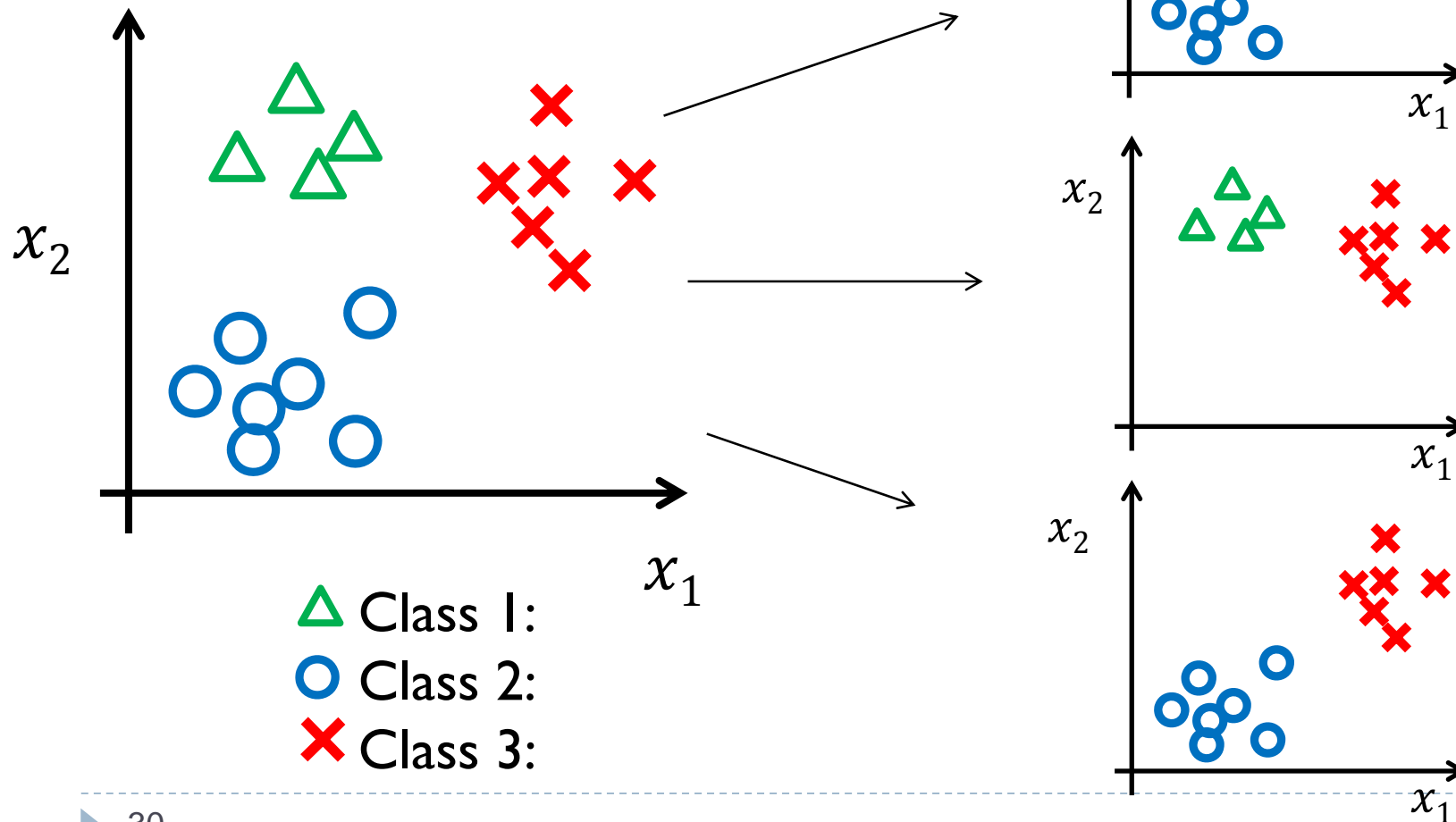
Multi-Category Classification

► One-vs-all (one-vs-rest)

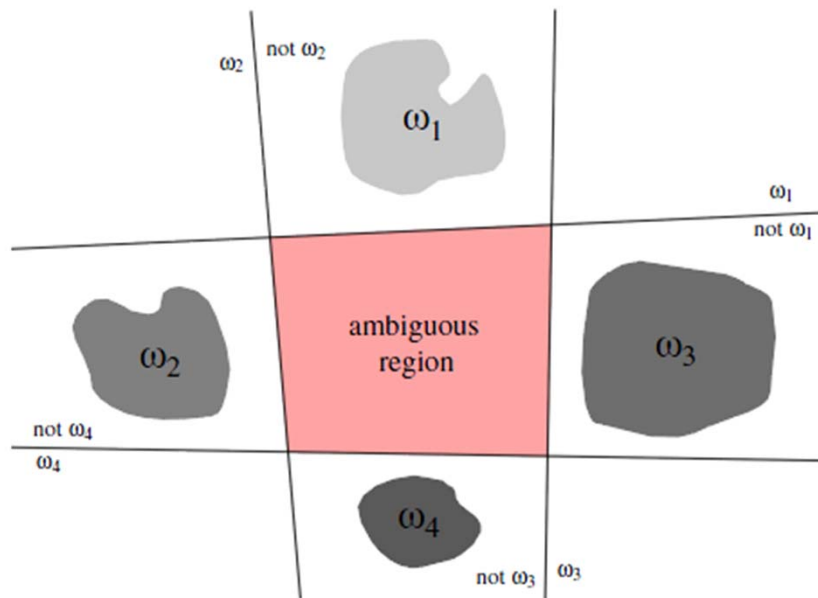


Multi-Category Classification

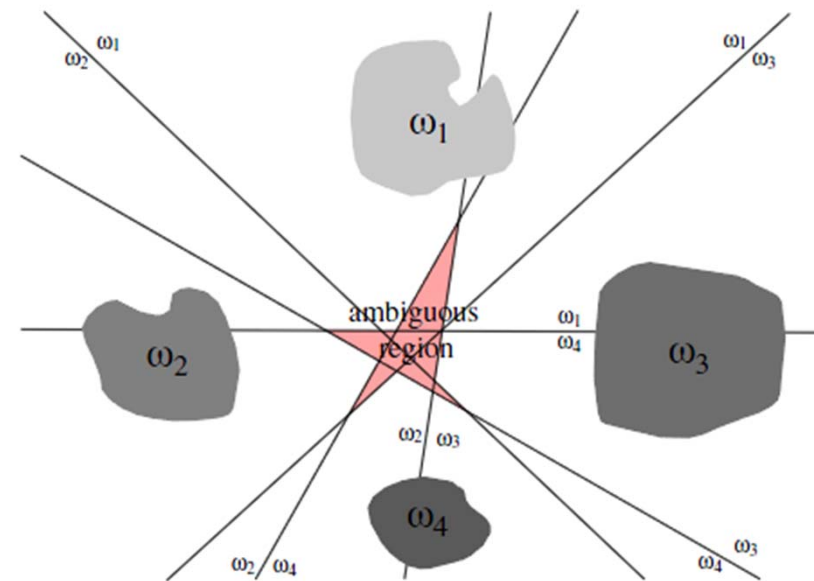
► One-vs-one



Multi-Category Classification: Ambiguity



one versus rest



one versus one

[Duda,Hart&Stork]

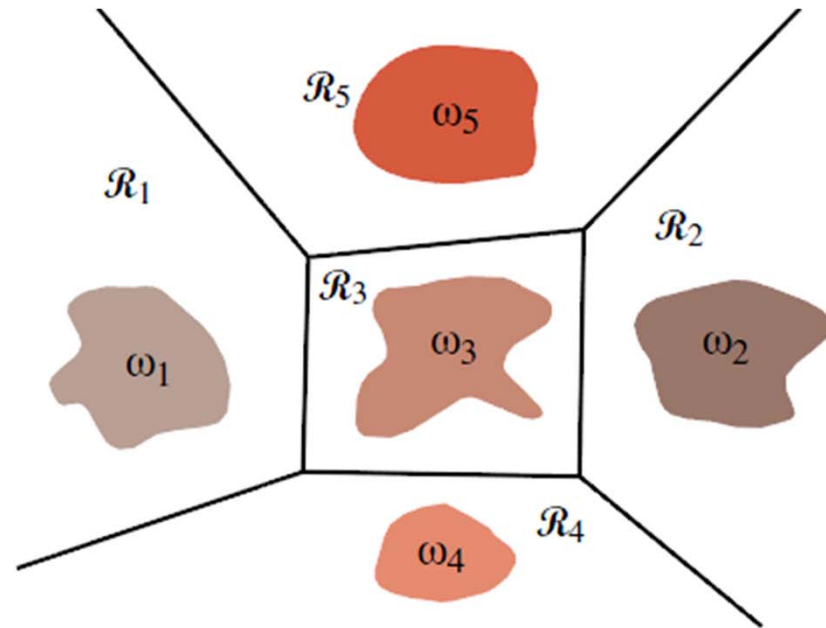
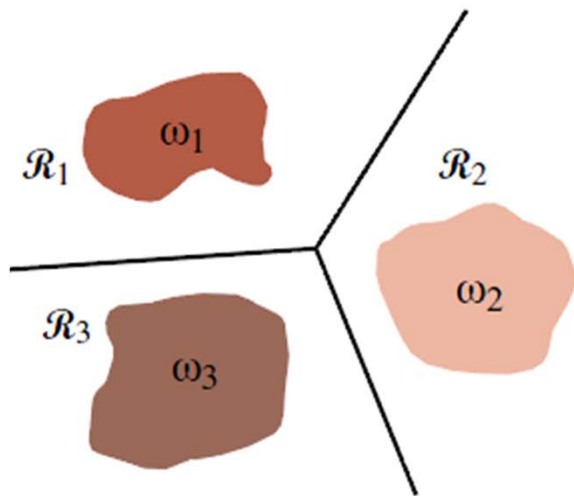
Multi-Category Classification: Linear Machine

- ▶ A discriminant function $g_i(\mathbf{x}) = \mathbf{w}_i^T \mathbf{x} + w_{i0}$ for each class ω_i ($i = 1, \dots, c$):
 - ▶ \mathbf{x} is assigned to class ω_i if:

$$g_i(\mathbf{x}) > g_j(\mathbf{x}) \quad \forall j \neq i$$

- ▶ **Decision surfaces** (boundaries) can also be found using discriminant functions
 - ▶ Boundary of the contiguous \mathcal{R}_i and \mathcal{R}_j : $\forall \mathbf{x}, g_i(\mathbf{x}) = g_j(\mathbf{x})$
 - ▶ $(\mathbf{w}_i - \mathbf{w}_j)^T \mathbf{x} + (w_{i0} - w_{j0}) = 0$

Multi-Category Classification: Linear Machine



[Duda,Hart&Stork]

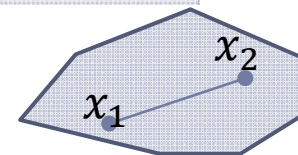
Multi-Category Classification: Linear Machine

- ▶ Decision regions are convex
 - ▶ Linear machines are most suitable for problems where $p(\mathbf{x}|\omega_i)$ are unimodal.

$$x_1, x_2 \in \mathcal{R}_i \Rightarrow \forall j \neq i, g_i(x_1) \geq g_j(x_1) \\ g_i(x_2) \geq g_j(x_2)$$

$$\Rightarrow \alpha g_i(x_1) + (1 - \alpha)g_i(x_2) \geq \alpha g_j(x_1) + (1 - \alpha)g_j(x_2) \\ g_i \text{ is linear} \Rightarrow g_i(\alpha x_1 + (1 - \alpha)x_2) \geq g_j(\alpha x_1 + (1 - \alpha)x_2) \\ \Rightarrow \alpha x_1 + (1 - \alpha)x_2 \in \mathcal{R}_i$$

Convex region definition: $\forall x_1, x_2 \in \mathcal{R}, 0 \leq \alpha \leq 1 \Rightarrow \alpha x_1 + (1 - \alpha)x_2 \in \mathcal{R}$



Multi-Category Classification: Target Coding Scheme

- ▶ Target values:

- ▶ Binary classification: a target variable $y \in \{0,1\}$

- ▶ Multiple classes ($K > 2$):

- ▶ TargetClass C_j : $y_j = 1$
 $\forall i \neq j \ y_i = 0$

SSE Cost Function: Multi-Class

$$J(\mathbf{W}) = \text{Tr}\{(X\mathbf{W} - \mathbf{Y})^T (X\mathbf{W} - \mathbf{Y})\}$$

???

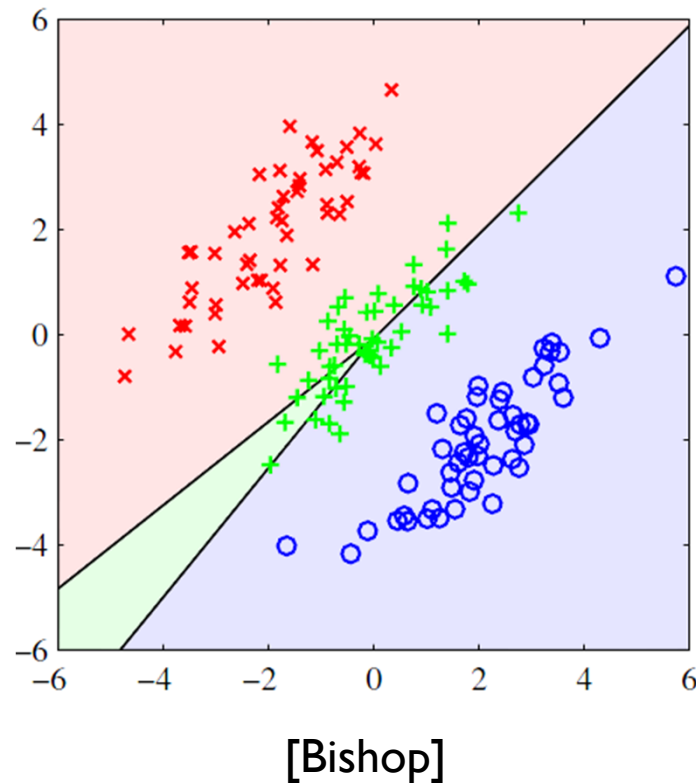
$$\mathbf{X} = \begin{bmatrix} 1 & x_1^{(1)} & \cdots & x_d^{(1)} \\ 1 & x_1^{(2)} & \cdots & x_d^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(n)} & \cdots & x_d^{(n)} \end{bmatrix} \quad \mathbf{W} = [\mathbf{w}_1 \quad \cdots \quad \mathbf{w}_c]$$

$$\mathbf{Y} = [\mathbf{y}_1 \quad \cdots \quad \mathbf{y}_c]$$

$$\nabla_{\mathbf{W}} J(\mathbf{W}) = \mathbf{0} \Rightarrow \widehat{\mathbf{W}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X} \mathbf{Y}$$

SSE Cost Function: Multi-Class

- Low performance of the SSE cost function for the classification problem



Perceptron: Multi-Class

$$\hat{y} = \operatorname{argmax}_{i=1,\dots,c} \mathbf{w}_i^T \mathbf{x}$$

$$J_P(\mathbf{W}) = - \sum_{i \in \mathcal{M}} \left(\mathbf{w}_{y^{(i)}} - \mathbf{w}_{\hat{y}^{(i)}} \right)^T \mathbf{x}^{(i)}$$

\mathcal{M} : subset of training data that are misclassified

$$\mathcal{M} = \{i | \hat{y}^{(i)} \neq y^{(i)}\}$$

Initialize $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_c]$, $k \leftarrow 0$

repeat

$k \leftarrow (k + 1) \bmod N$

 if $\mathbf{x}^{(i)}$ is misclassified then

$$\mathbf{w}_{\hat{y}^{(i)}} = \mathbf{w}_{\hat{y}^{(i)}} - \mathbf{x}^{(i)}$$

$$\mathbf{w}_{y^{(i)}} = \mathbf{w}_{y^{(i)}} + \mathbf{x}^{(i)}$$

Until all patterns properly classified

Generalized LDFs

- ▶ Linear combination of a fixed non-linear functions of the input vector

$$g(\mathbf{x}; \mathbf{w}) = w_0 + w_1 \phi_1(\mathbf{x}) + \dots w_m \phi_m(\mathbf{x})$$

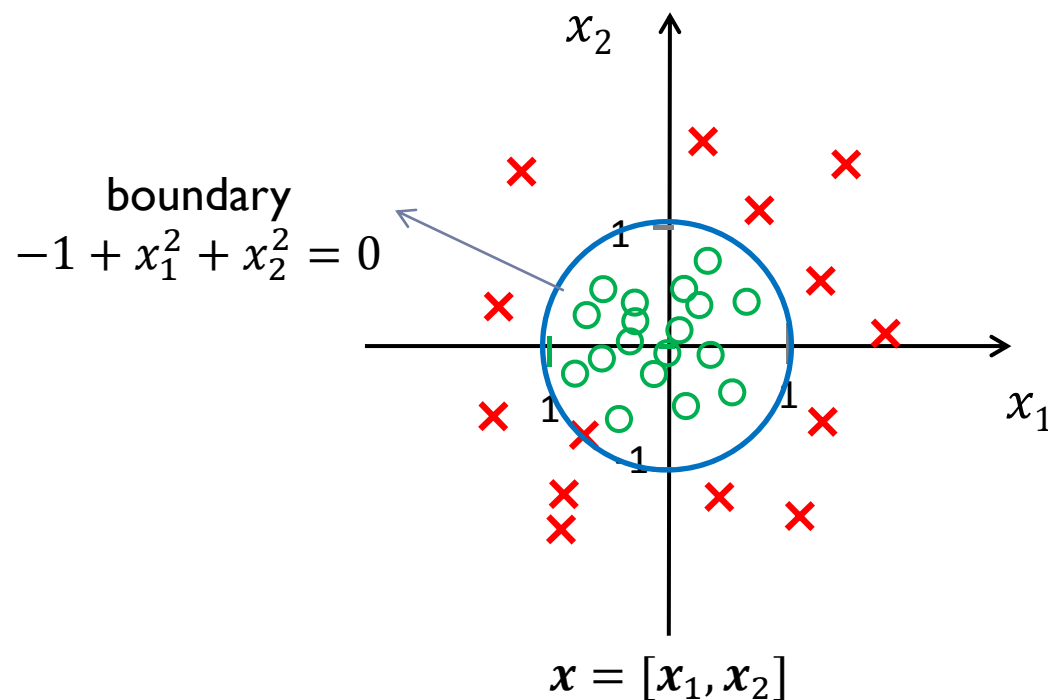
/phi

$\{\phi_1(\mathbf{x}), \dots, \phi_m(\mathbf{x})\}$: set of basis functions (or features)

$$\phi_i(\mathbf{x}): \mathbb{R}^d \rightarrow \mathbb{R}$$

Generalized LDFs: Example

- ▶ Choose non-linear features
- ▶ Discriminant functions are still linear in parameters \mathbf{w}



$$\phi(\mathbf{x}) = [1, x_1, x_2, x_1^2, x_2^2, x_1x_2]$$

$$\mathbf{w} = [-1, 0, 0, 1, 1, 0]$$

if $\mathbf{w}^T \phi(\mathbf{x}) \geq 0$ then $y = 1$
else $y = -1$