

Single-Layer Perceptron Classifiers

Berlin Chen, 2002

Outline

- Foundations of trainable decision-making networks to be formulated
 - Input space to output space (classification space)
- Focus on the classification of linearly separable classes of patterns
 - Linear discriminating functions and simple correction function
 - Continuous error function minimization
- Explanation and justification of perceptron and delta training rules

Classification Model, Features, and Decision Regions

- A pattern is the quantitative description of an object, event, or phenomenon
 - **Spatial patterns**: weather maps, fingerprints ...
 - **Temporal patterns**: speech signals ...
- Pattern classification/recognition
 - Assign the input data (a physical object, event, or phenomenon) to one of the pre-specified classes (categories)
 - Discriminate the input data within object population via the search for **invariant attributes** among members of the population

Classification Model, Features, and Decision Regions (cont.)

- The block diagram of the recognition and classification system

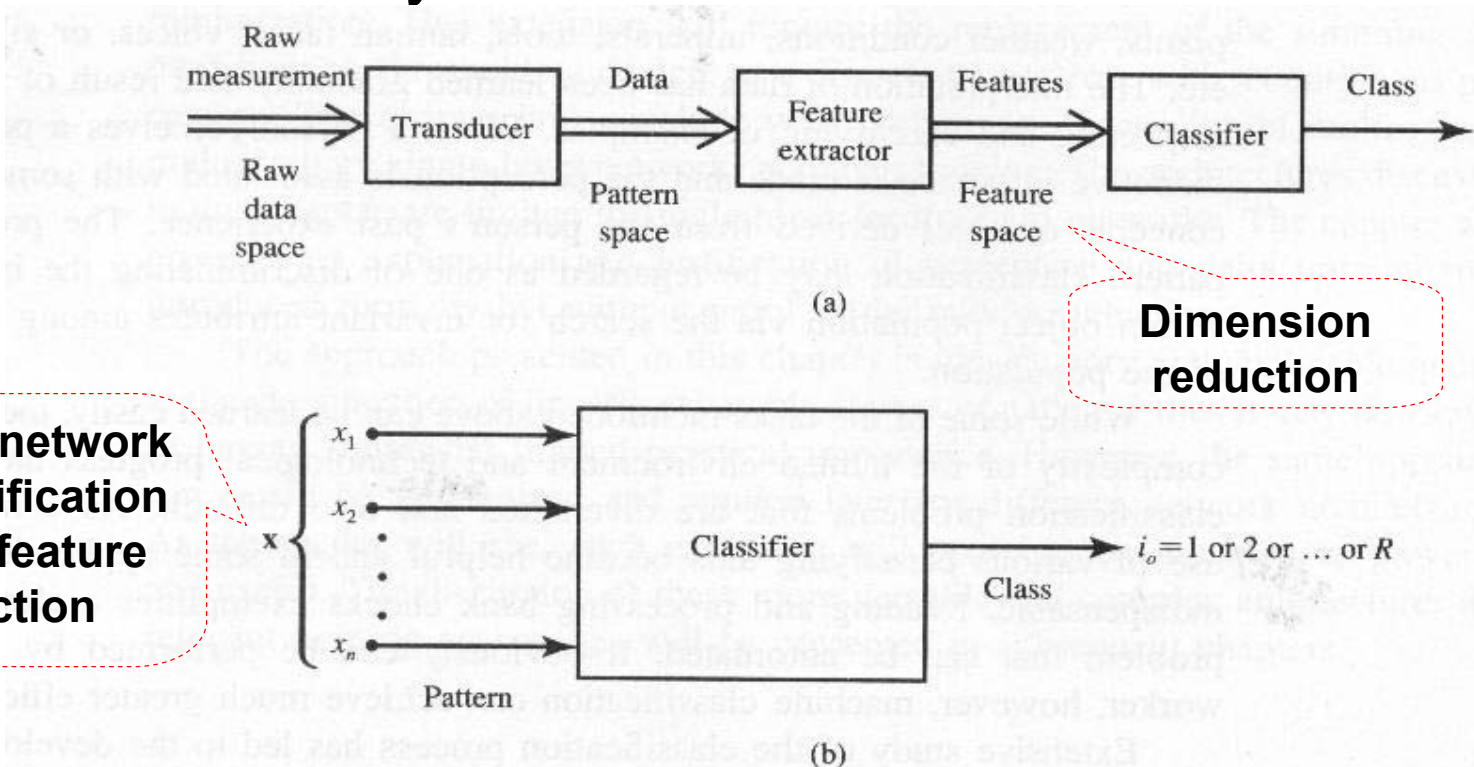


Figure 3.1 Recognition and classification system: (a) overall block diagram and (b) pattern classifier.

Classification Model, Features, and Decision Regions (cont.)

- More about Feature Extraction
 - The compressed data from the input patterns while poses salient information
 - E.g.
 - Speech vowel sounds analyzed in 16-channel filterbanks can provide 16 spectral vectors, which can be further transformed into two dimensions
 - Tone height (high-low) and retraction (front-back)
 - Input patterns to be projected and reduced to lower dimensions

Classification Model, Features, and Decision Regions (cont.)

- More about Feature Extraction

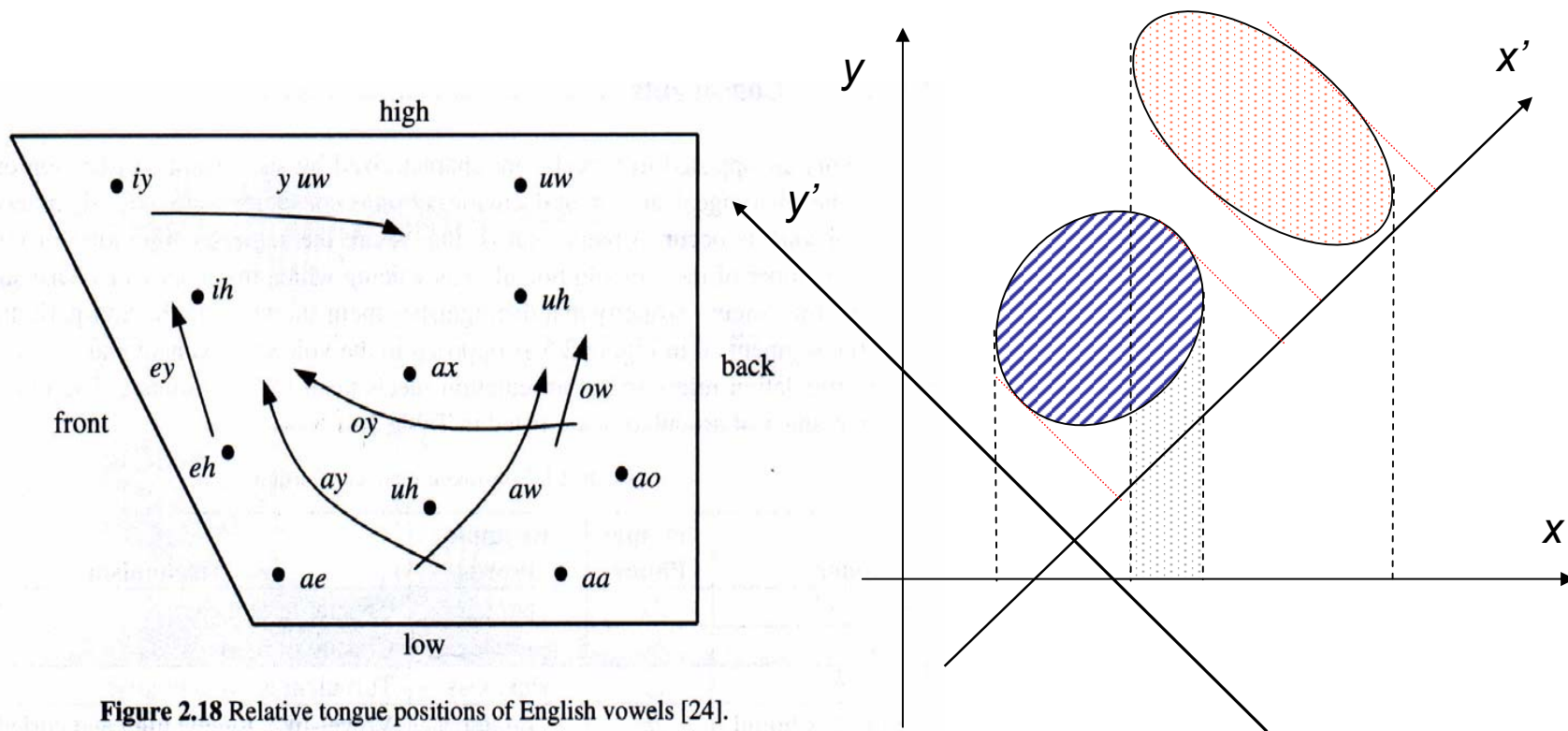


Figure 2.18 Relative tongue positions of English vowels [24].

Classification Model, Features, and Decision Regions (cont.)

- Two simple ways to generate the **pattern vectors** for cases of **spatial** and **temporal** objects to be classified

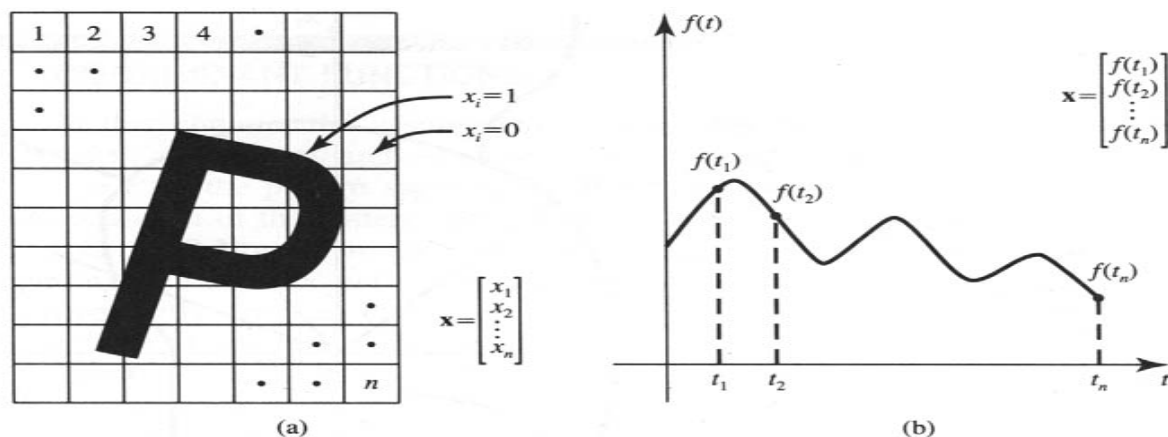


Figure 3.2 Two simple ways of coding patterns into pattern vectors: (a) spatial object and (b) temporal object (waveform).

- A pattern classifier maps input patterns (vectors) in E^n space into numbers (E^1) which specify the membership

$$j = i_0(\mathbf{x}), \quad j = 1, 2, \dots, R$$

Classification Model, Features, and Decision Regions (cont.)

- Classification described in geometric terms

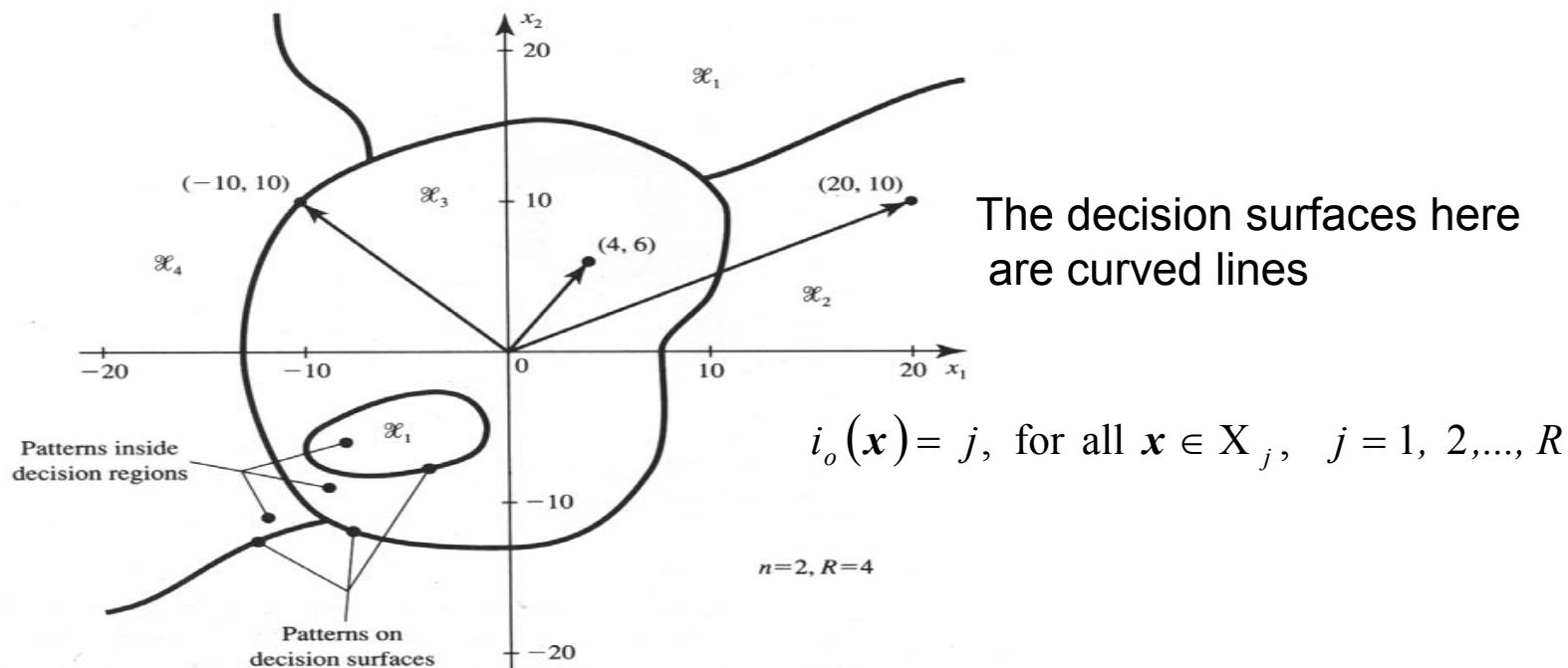


Figure 3.3 Decision regions example.

- **Decision regions**
- **Decision surfaces:** generally, the decision surfaces for n -dimensional patterns may be $(n-1)$ -dimensional hyper-surfaces

Discriminant Functions

- Determine the membership in a category by the classifier based on the comparison of R **discriminant functions** $g_1(\mathbf{x}), g_2(\mathbf{x}), \dots, g_R(\mathbf{x})$
 - When \mathbf{x} is within the region \mathbf{X}_k if $g_k(\mathbf{x})$ has the largest value $i_0(\mathbf{x}) = k$ if $g_k(\mathbf{x}) > g_j(\mathbf{x})$ for $k, j = 1, 2, \dots, R, k \neq j$

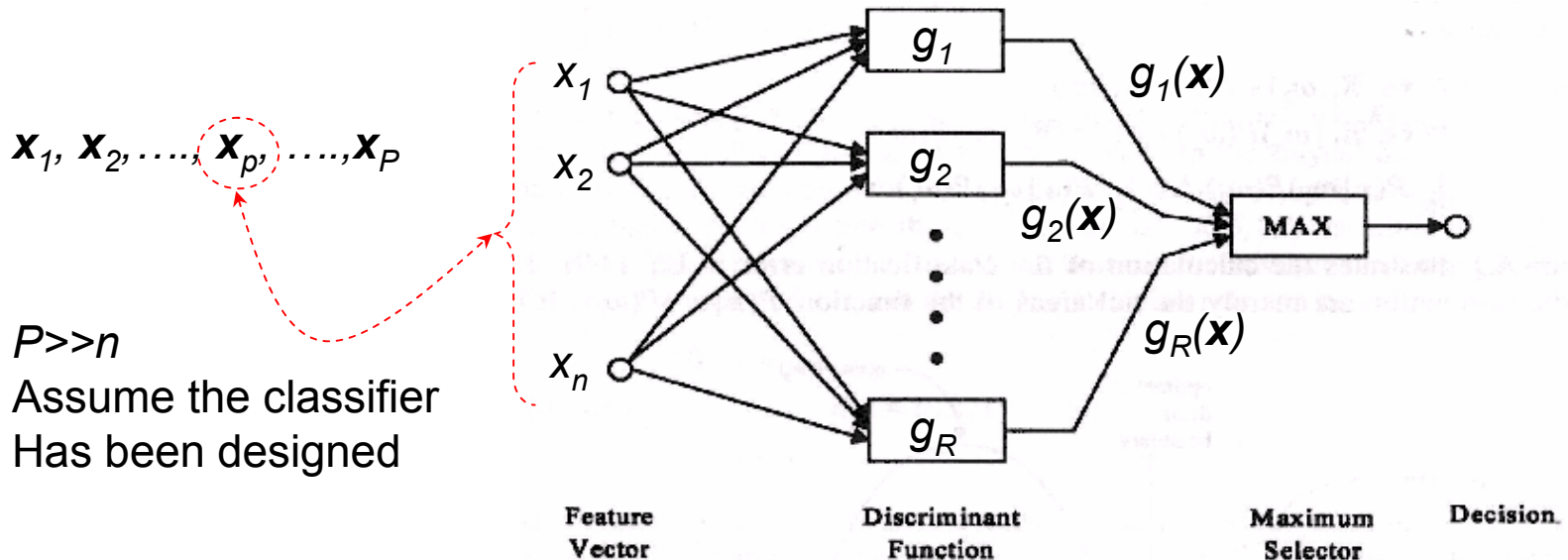


Figure 4.2 Block diagram of a classifier based on discriminant functions [22].

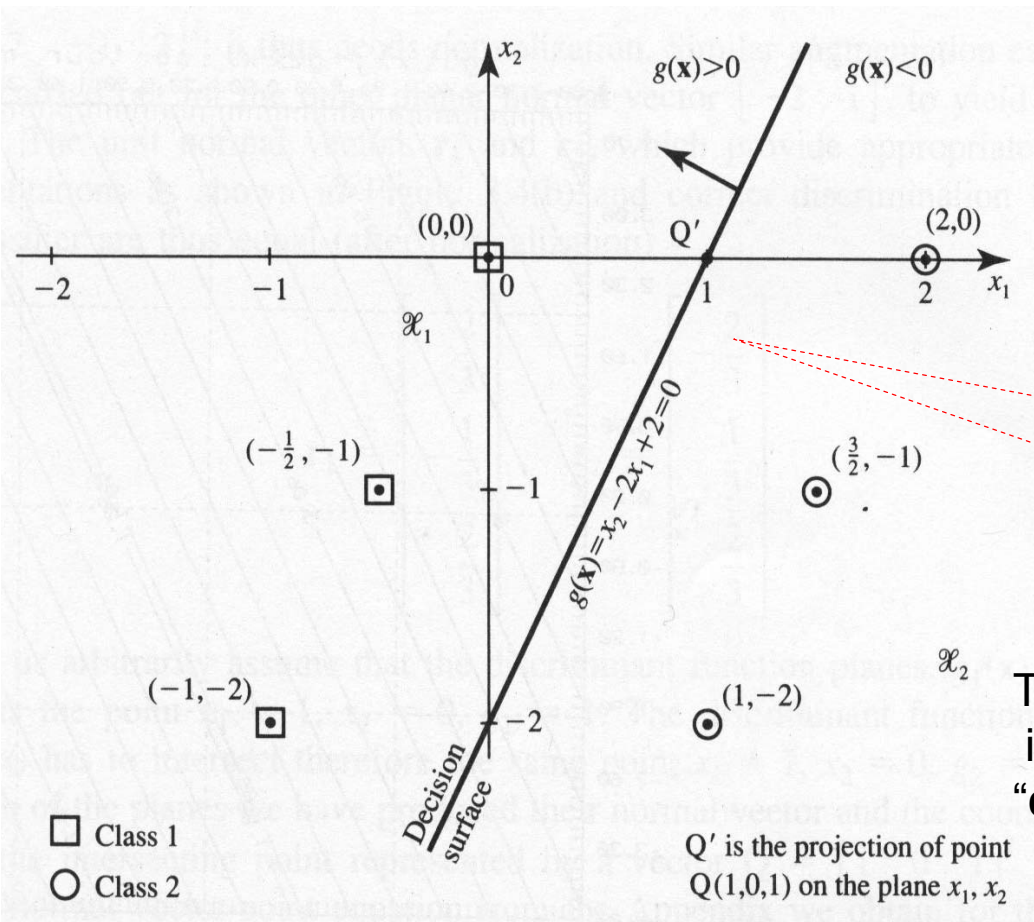
Discriminant Functions (cont.)

- Example 3.1

Decision surface Equation: $g(\mathbf{x}) = g_1(\mathbf{x}) - g_2(\mathbf{x})$
 $= -2x_1 + x_2 + 2$

$g(\mathbf{x}) > 0$: class1

$g(\mathbf{x}) < 0$: class2



The **decision surface** does not uniquely specify the discriminant functions

The classifier that classifies patterns into two classes or categories is called “**dichotomizer**”

“two” “cut”

Discriminant Functions (cont.)

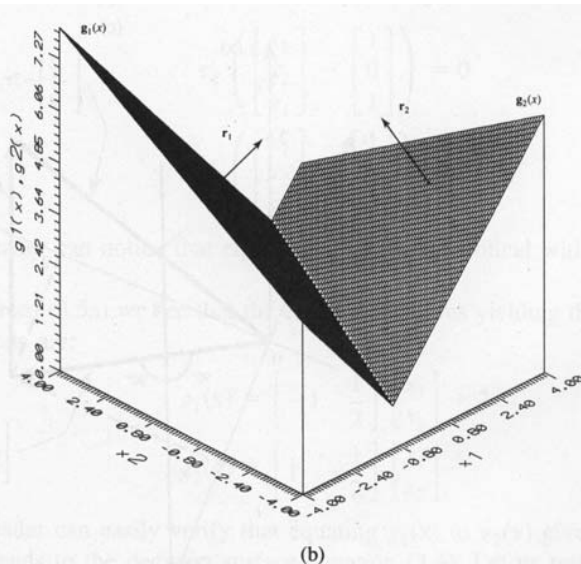


Figure 3.4a,b Illustration for Example 3.1: (a) pattern display and decision surface, (b) discriminant functions.

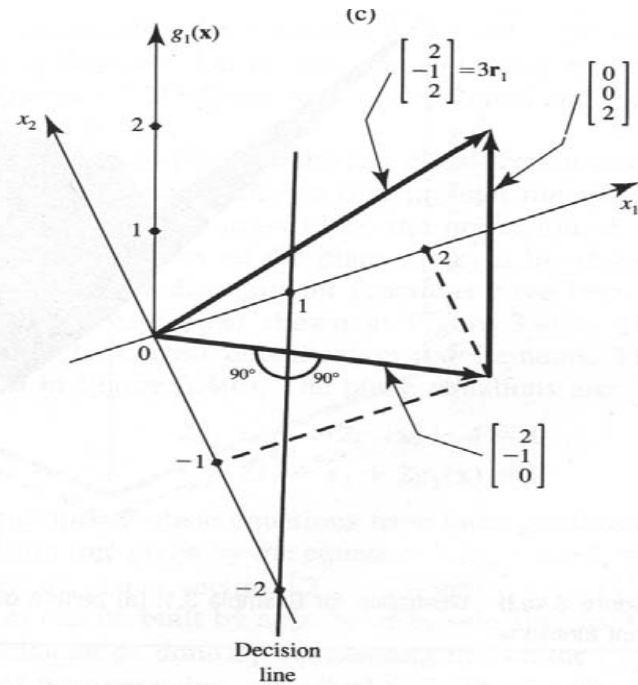


Figure 3.4c,d Illustration for Example 3.1 (continued): (c) contour map of discriminant functions, and (d) construction of the normal vector for $a_1(x)$.

Discriminant Functions (cont.)

$$(x-0, y+2, g_1-1)(2, -1, 1)=0$$

$$2x-y-2+g_1-1=0$$

$$g_1=-2x+y+3$$

$$(x-0, y+2, g_2-1)(-2, 1, 1)=0$$

$$-2x+y+2+g_2-1=0$$

$$g_2=2x-y-1$$

$$g=g_1-g_2=0$$

$$-4x+2y+4=0$$

$$-2x+y+2=0$$

Solution 1

$$g_1(\mathbf{x}) = \begin{bmatrix} -2 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + 3$$

$$g_2(\mathbf{x}) = \begin{bmatrix} 2 & -1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$(x-0, y+2, g_1-1)(2, -1, 2)=0$$

$$2x-y-2+2g_1-2=0$$

$$g_1=-x+1/2y+2$$

$$(x-0, y+2, g_2-1)(-2, 1, 2)=0$$

$$-2x+y+2+2g_2-2=0$$

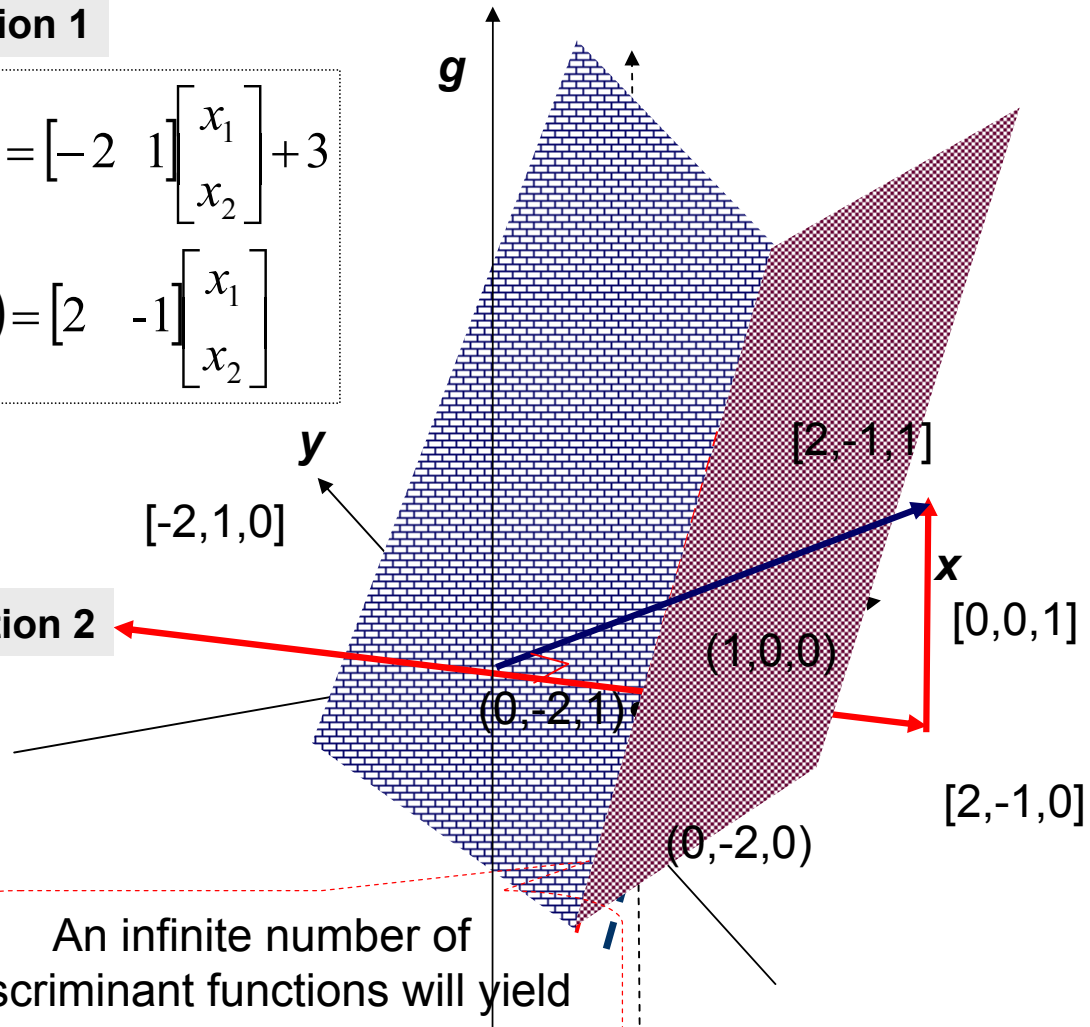
$$g_2=x-1/2y$$

$$g=g_1-g_2=0$$

$$-2x+y+2=0$$

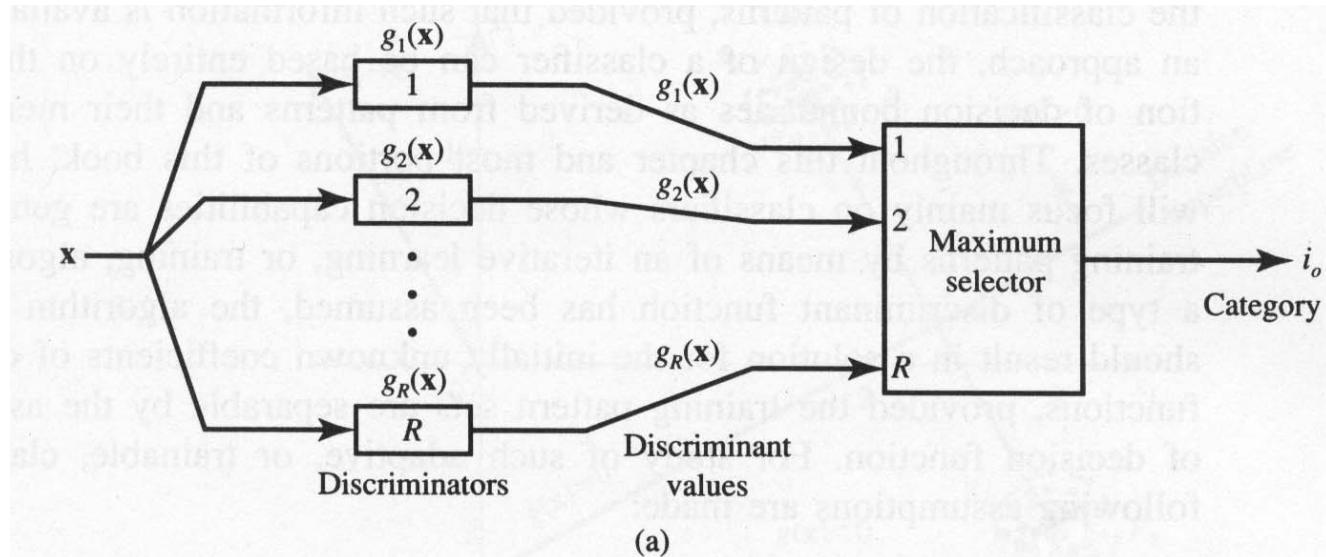
Solution 2

An infinite number of discriminant functions will yield correct classification

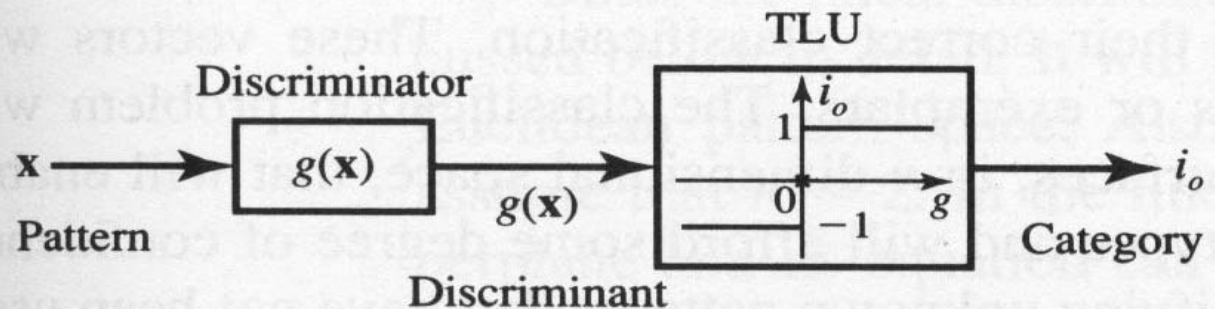


Discriminant Functions (cont.)

Multi-class



Two-class



$$g(\mathbf{x}) = g_1(\mathbf{x}) - g_2(\mathbf{x})$$

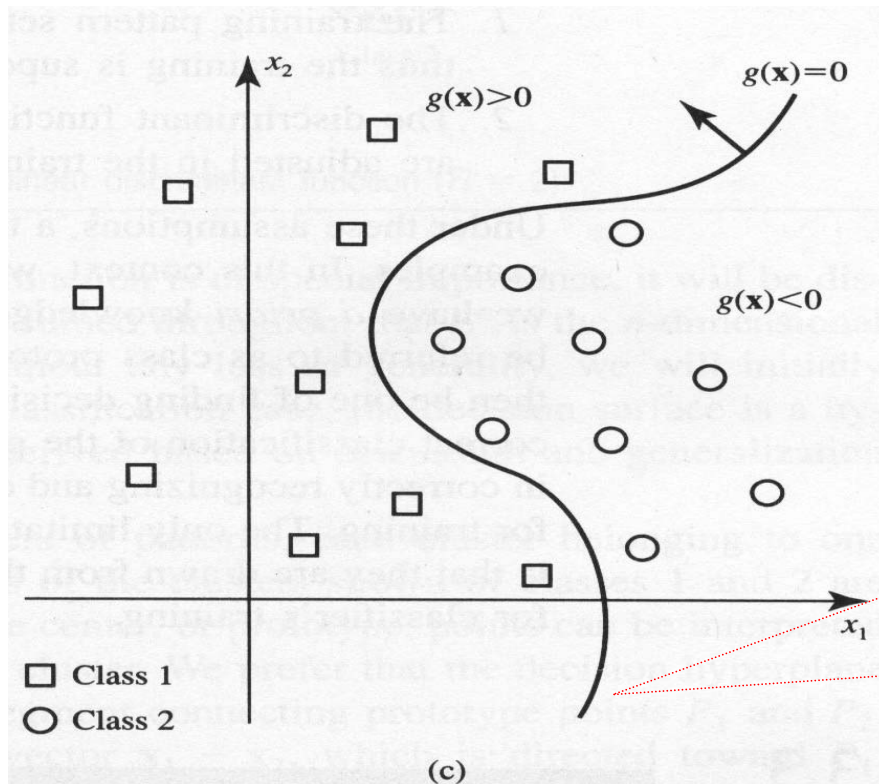
$$g(\mathbf{x}) > 0 : \text{class 1}$$

$$g(\mathbf{x}) < 0 : \text{class 2}$$

subtraction

Sign examination

Discriminant Functions (cont.)



(a) into R categories, (b) dichotomizer ($R = 2$), and

The design of discriminator for this case is not straightforward. The discriminant functions may result as nonlinear functions of x_1 and x_2

Bayes' Decision Theory

- A decision-making based on both the posterior knowledge obtained from specific observation data and prior knowledge of the categories
 - Prior class probabilities $P(\omega_i), \forall \text{ class } i$
 - Class-conditioned probabilities $P(x|\omega_i), \forall \text{ class } i$

$$k = \arg \max_i P(\omega_i|x) = \arg \max_i \frac{P(x|\omega_i)P(\omega_i)}{P(x)} = \arg \max_i \frac{P(x|\omega_i)P(\omega_i)}{\sum_{j=1} P(x|\omega_j)P(\omega_j)}$$

$$k = \arg \max_i P(\omega_i|x) = \arg \max_i P(x|\omega_i)P(\omega_i)$$

Bayes' Decision Theory (cont.)

- Bayes' decision rule designed to minimize the overall risk involved in making decision

- The expected loss (**conditional risk**) when making decision δ_i

$$\begin{aligned} R(\delta_i|x) &= \sum_j l(\delta_i|\omega_j, x)P(\omega_j|x), \quad \text{where} \quad l(\delta_i|\omega_j, x) = \begin{cases} 0, & i = j \\ 1, & i \neq j \end{cases} \\ &= \sum_{j \neq i} P(\omega_j|x) \\ &= 1 - P(\omega_i|x) \end{aligned}$$

- The overall risk (**Bayes' risk**)

$$R = \int_{-\infty}^{\infty} R(\delta(x)|x)p(x)dx, \quad \delta(x): \text{the selected decision for a sample } x$$

- **Minimize the overall risk** (classification error) by computing the conditional risks and select the decision δ_i for which the conditional risk $R(\delta_i|x)$ is minimum, i.e., $P(\omega_i|x)$ is maximum (**minimum-error-rate decision rule**)

Bayes' Decision Theory (cont.)

- Two-class pattern classification

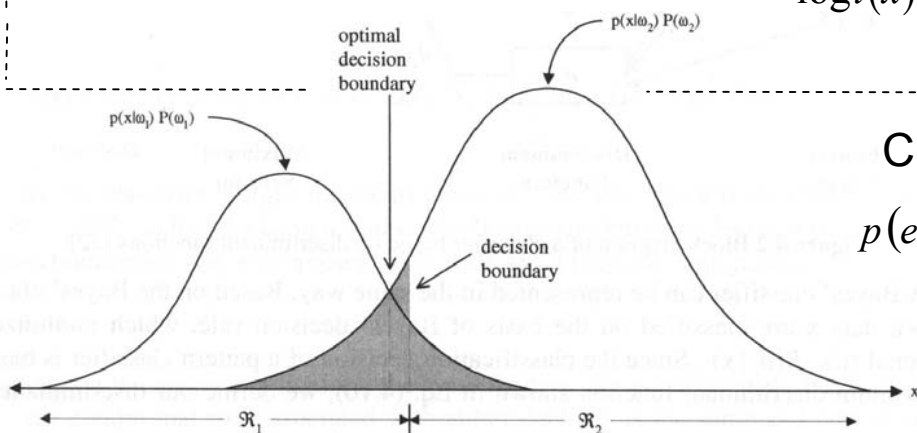
$$g_1(x) = P(\omega_1|x) \cong P(x|\omega_1)P(\omega_1), \quad g_2(x) = P(\omega_2|x) \cong P(x|\omega_2)P(\omega_2)$$

Bayes' Classifier

Likelihood ratio or log-likelihood ratio:

$$P(x|\omega_1)P(\omega_1) \underset{\omega_2}{>} P(x|\omega_2)P(\omega_2) \iff l(x) = \frac{P(x|\omega_1)}{P(x|\omega_2)} \underset{\omega_2}{>} \frac{P(\omega_2)}{P(\omega_1)}$$

$$\log l(x) = \log P(x|\omega_1) - \log P(x|\omega_2) \underset{\omega_2}{>} \log P(\omega_2) - \log P(\omega_1)$$



Classification error:

$$\begin{aligned} p(\text{error}) &= P(x \in R_1, \omega_2) + P(x \in R_2, \omega_1) \\ &= P(x \in R_1 | \omega_2)P(\omega_2) + P(x \in R_2 | \omega_1)P(\omega_1) \\ &= \int_{R_1} P(x|\omega_2)P(\omega_2)dx + \int_{R_2} P(x|\omega_1)P(\omega_1)dx \end{aligned}$$

Figure 4.1 Calculation of the likelihood of classification error [22]. The shaded area represents the integral value in Eq. (4.9).

Bayes' Decision Theory (cont.)

- When the environment is multivariate Gaussian, the Bayes' classifier reduces to a linear classifier
 - The same form taken by the perceptron
 - But the linear nature of the perceptron is not contingent on the assumption of Gaussianity

$$P(\mathbf{x}|\omega) = \frac{1}{(2\pi)^{n/2} |\boldsymbol{\Sigma}|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^t \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right]$$

$$\text{Class } \omega_1 : E[\mathbf{X}] = \boldsymbol{\mu}_1$$

$$E[(\mathbf{X} - \boldsymbol{\mu}_1)(\mathbf{X} - \boldsymbol{\mu}_1)^t] = \boldsymbol{\Sigma}$$

$$P(\omega_1) = P(\omega_2) = \frac{1}{2}$$

$$\text{Class } \omega_2 : E[\mathbf{X}] = \boldsymbol{\mu}_2$$

$$E[(\mathbf{X} - \boldsymbol{\mu}_2)(\mathbf{X} - \boldsymbol{\mu}_2)^t] = \boldsymbol{\Sigma}$$

Assumptions

Bayes' Decision Theory (cont.)

- When the environment is Gaussian, the Bayes' classifier reduces to a linear classifier (cont.)

$$\begin{aligned}\log l(\mathbf{x}) &= \log P(\mathbf{x}|\omega_1) - \log P(\mathbf{x}|\omega_2) \\ &= -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_1)^t \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}_1) + \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_2)^t \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}_2) \\ &= (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^t \boldsymbol{\Sigma}^{-1} \mathbf{x} + \frac{1}{2}(\boldsymbol{\mu}_2^t \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_2 - \boldsymbol{\mu}_1^t \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1) \\ &= \mathbf{w}\mathbf{x} + b\end{aligned}$$

$$\therefore \log l(\mathbf{x}) = \mathbf{w}\mathbf{x} + b \begin{matrix} > \\ < \end{matrix} \begin{matrix} \omega_1 \\ \omega_2 \end{matrix}$$

Bayes' Decision Theory (cont.)

- Multi-class pattern classification

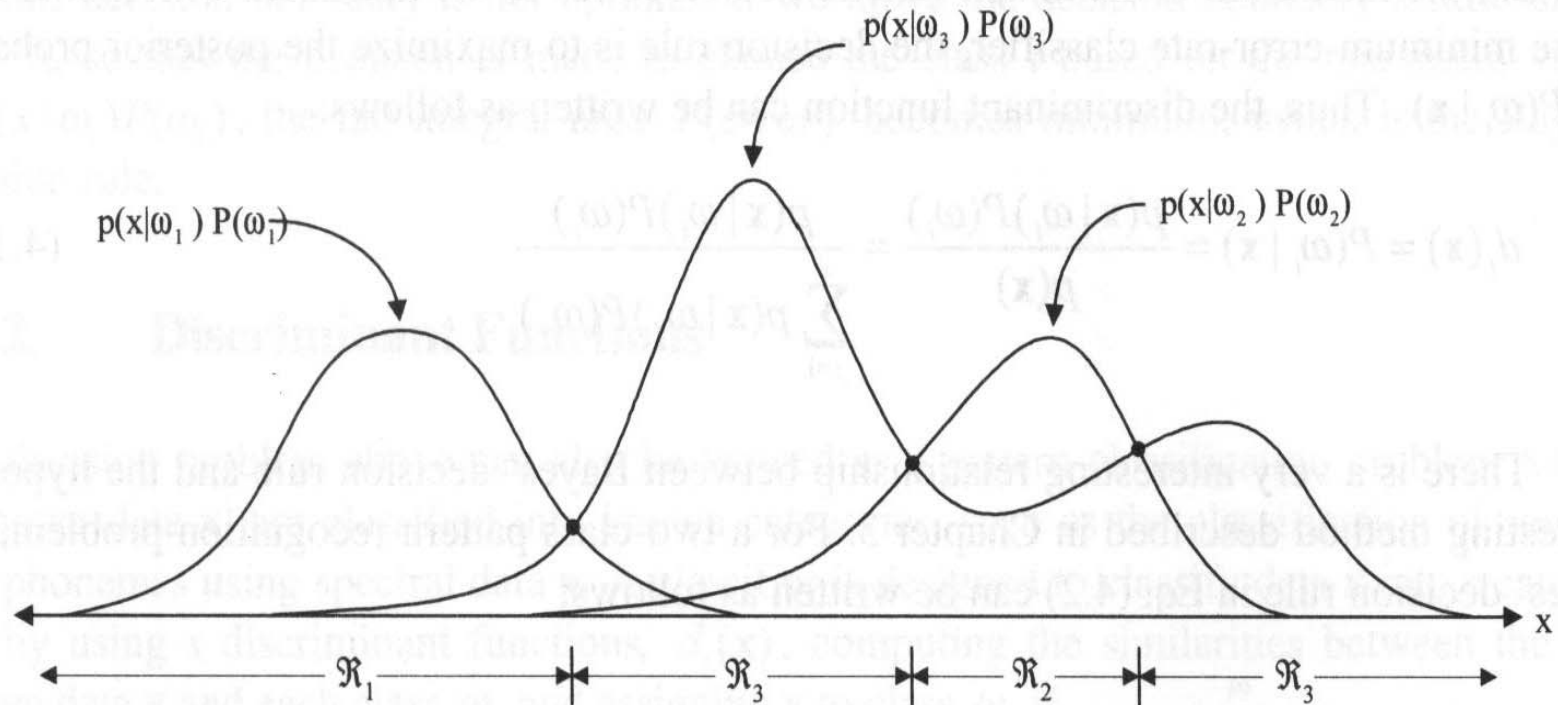


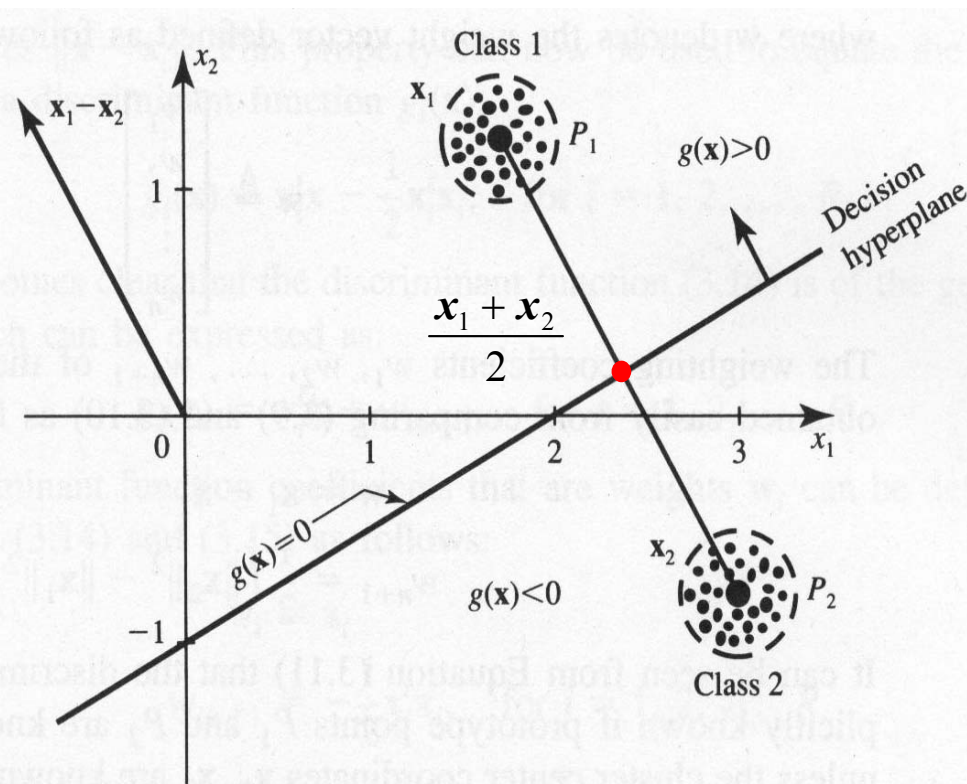
Figure 4.3 An example of decision boundaries and regions. For simplicity, we use scalar variable x instead of a multi-dimensional vector [22].

Linear Machine and Minimum Distance Classification

- Find the linear-form discriminant function for two-class classification when the class prototypes are known
- **Example 3.1:** Select the decision hyperplane that contains the midpoint of the line segment connecting center point of two classes

Linear Machine and Minimum Distance Classification (cont.)

The dichotomizer's discriminant function $g(\mathbf{x})$:



$$(\mathbf{x}_1 - \mathbf{x}_2)^t \left(\mathbf{x} - \frac{\mathbf{x}_1 + \mathbf{x}_2}{2} \right) = 0$$

$$(\mathbf{x}_1 - \mathbf{x}_2)^t \mathbf{x} + \frac{1}{2} (\|\mathbf{x}_2\|^2 - \|\mathbf{x}_1\|^2) = 0$$

Taken as $\begin{bmatrix} \mathbf{w} \\ w_{n+1} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix} = 0$, where

$$\mathbf{w} = \mathbf{x}_1 - \mathbf{x}_2$$

$$w_{n+1} = \frac{1}{2} (\|\mathbf{x}_2\|^2 - \|\mathbf{x}_1\|^2)$$

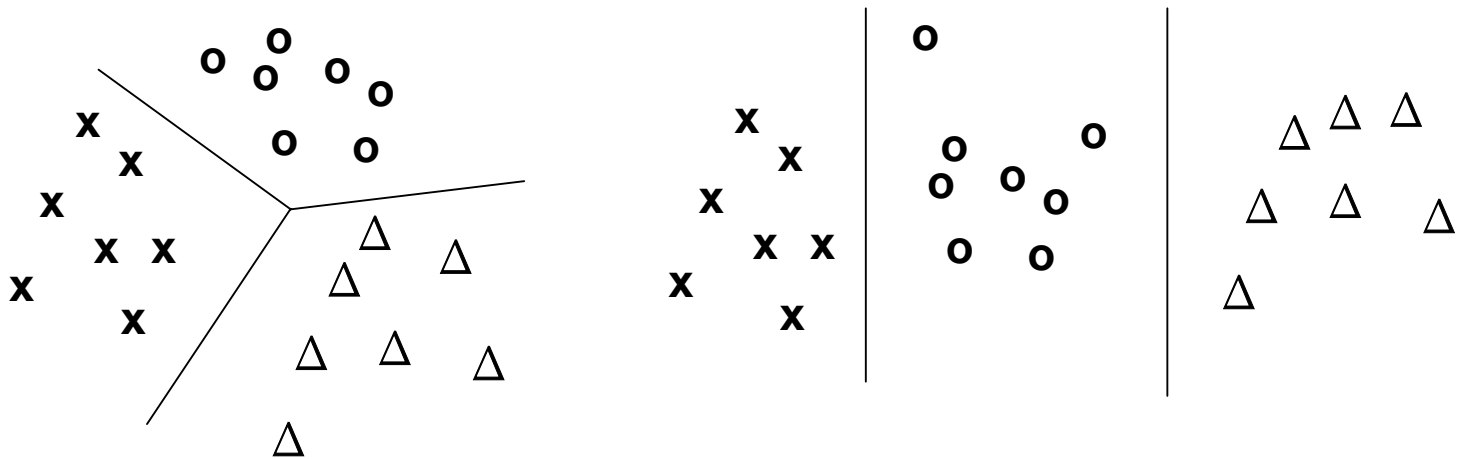
Augmented input pattern

It is a simple minimum-distance classifier.

Linear Machine and Minimum Distance Classification (cont.)

- The linear-form discriminant functions for multi-class classification
 - There are up to $R(R-1)/2$ decision hyperplanes for R pairwise separable classes

Some classes may not be contiguous



Linear Machine and Minimum Distance Classification (cont.)

- Linear machine or minimum-distance classifier
 - Assume the class prototypes are known for all classes

- Euclidean distance between input pattern \mathbf{x} and the center of class i , \mathbf{x}_i :

$$\|\mathbf{x} - \mathbf{x}_i\| = \sqrt{(\mathbf{x} - \mathbf{x}_i)^t (\mathbf{x} - \mathbf{x}_i)}$$

- Minimizing $\|\mathbf{x} - \mathbf{x}_i\|^2 = \mathbf{x}^t \mathbf{x} - 2\mathbf{x}_i^t \mathbf{x} + \mathbf{x}_i^t \mathbf{x}_i$ is equal to

$$\text{maximizing } \mathbf{x}_i^t \mathbf{x} - \frac{1}{2} \mathbf{x}_i^t \mathbf{x}_i$$

The same for all classes

- Set the discriminant function for each class i to be:

$$g_i(\mathbf{x}) = \mathbf{x}_i^t \mathbf{x} - \frac{1}{2} \mathbf{x}_i^t \mathbf{x}_i \quad \longrightarrow \quad g_i(\mathbf{x}) = \mathbf{w}_i^t \mathbf{y}$$

$$g_i(\mathbf{x}) = \begin{bmatrix} \mathbf{w}_i \\ w_{i,n+1} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}, \text{ where}$$

$$\begin{aligned} \mathbf{w}_i &= \mathbf{x}_i \\ w_{i,n+1} &= -\frac{1}{2} (\mathbf{x}_i^t \mathbf{x}_i) \end{aligned}$$

Linear Machine and Minimum Distance Classification (cont.)

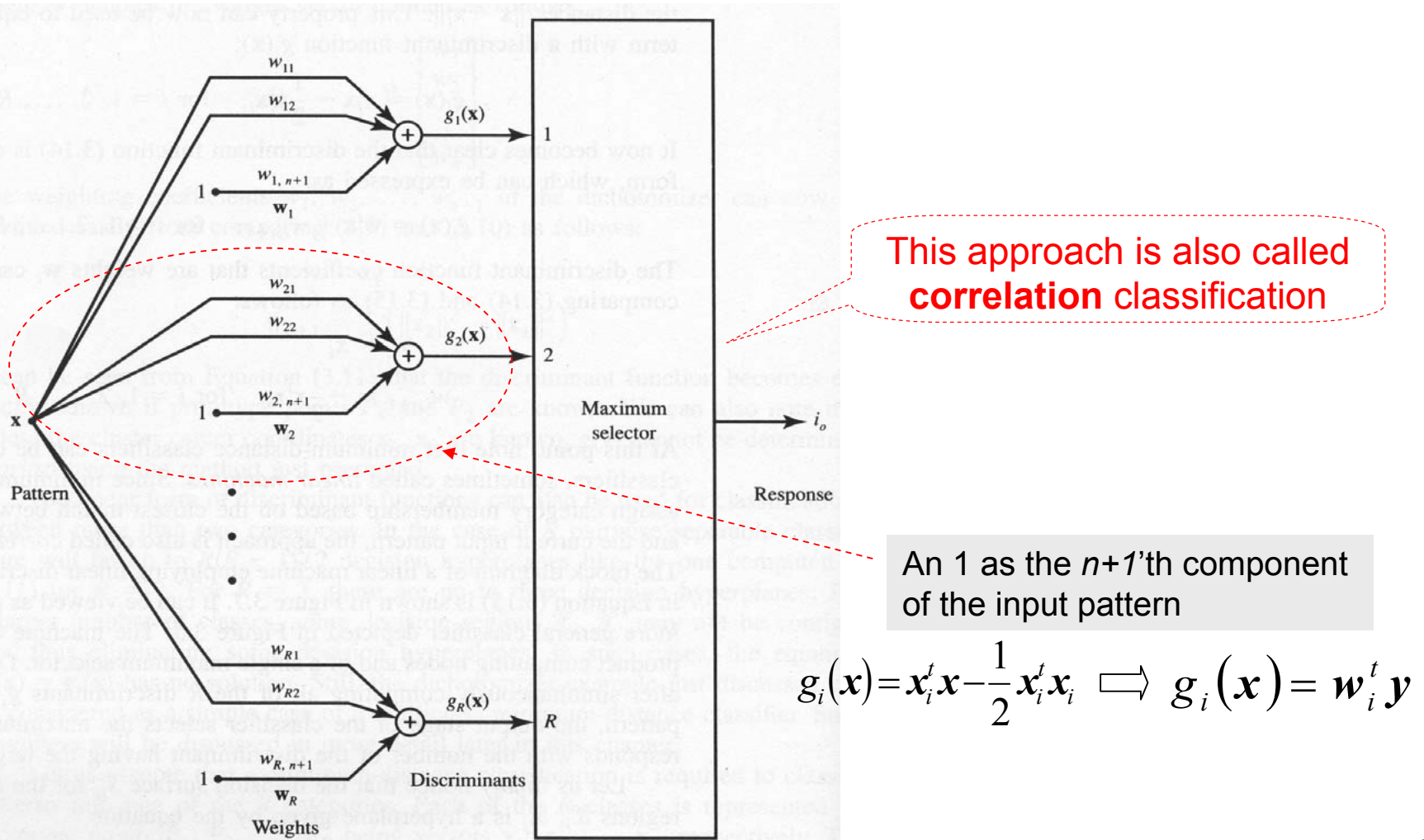


Figure 3.7 A linear classifier.

Linear Machine and Minimum Distance Classification (cont.)

- Example 3.2

$$w_1 = \begin{bmatrix} 10 \\ 2 \\ -52 \end{bmatrix}, \quad w_2 = \begin{bmatrix} 2 \\ -5 \\ -14.5 \end{bmatrix}, \quad w_3 = \begin{bmatrix} -5 \\ 5 \\ -25 \end{bmatrix}$$

$$g_1(x) = 10x_1 + 2x_2 - 52$$

$$g_2(x) = 2x_1 - 5x_2 - 14.5$$

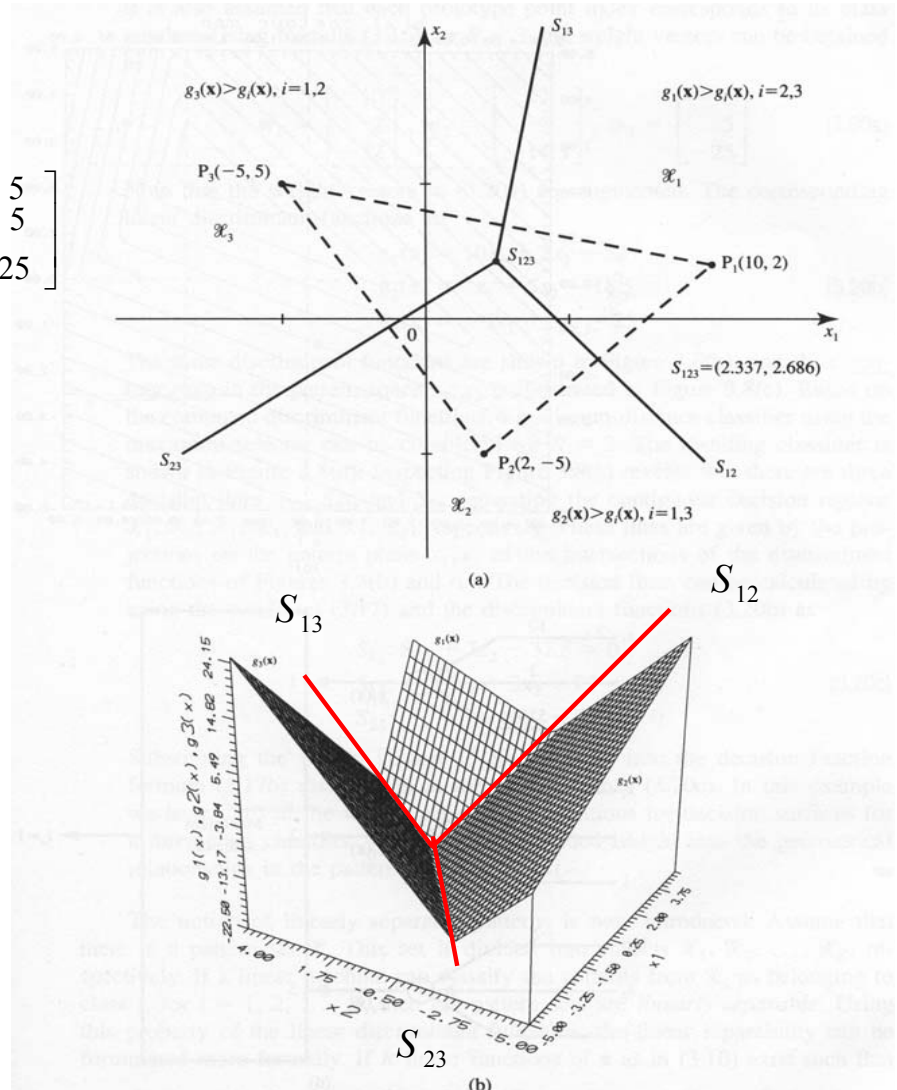
$$g_3(x) = -5x_1 + 5x_2 - 25$$

$$S_{12} : 8x_1 + 7x_2 - 37.5 = 0$$

$$S_{13} : -15x_1 + 3x_2 + 27 = 0$$

$$S_{23} : -7x_1 + 10x_2 - 10.5 = 0$$

$$g_i(x) = x_i^t x - \frac{1}{2} x_i^t x_i$$



Linear Machine and Minimum Distance Classification (cont.)

- If R linear discriminant functions exist for a set of patterns such that

$$g_i(\mathbf{x}) > g_j(\mathbf{x}) \text{ for } \mathbf{x} \in \text{Class } i,$$

$$i = 1, 2, \dots, R, j = 1, 2, \dots, R, i \neq j$$

- The classes are linearly separable

Linear Machine and Minimum Distance Classification (cont.)

P3.3 For the minimum-distance (linear) dichotomizer, the weight and augmented pattern vectors are

$$\mathbf{w} = \begin{bmatrix} 2 \\ -1 \\ 2 \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} x_1 \\ x_2 \\ 1 \end{bmatrix}$$

- (a) Find the equation of the decision surface in the pattern space.
- (b) Find the equation of the decision surface in the augmented pattern space.
- (c) Compute the new solution weight vector if the two class prototype points are

$$\mathbf{x}_1 = [2 \ 5]^t \text{ and } \mathbf{x}_2 = [-1 \ -3]^t.$$

- (d) Sketch the decision surfaces for each case in parts (a), (b), and (c).

Linear Machine and Minimum Distance Classification (cont.)

(a) $2x_1 - x_2 + 2 = 0$, decision surface is a line

(b) $2x_1 - x_2 + 2 = 0$, decision surface is a plane

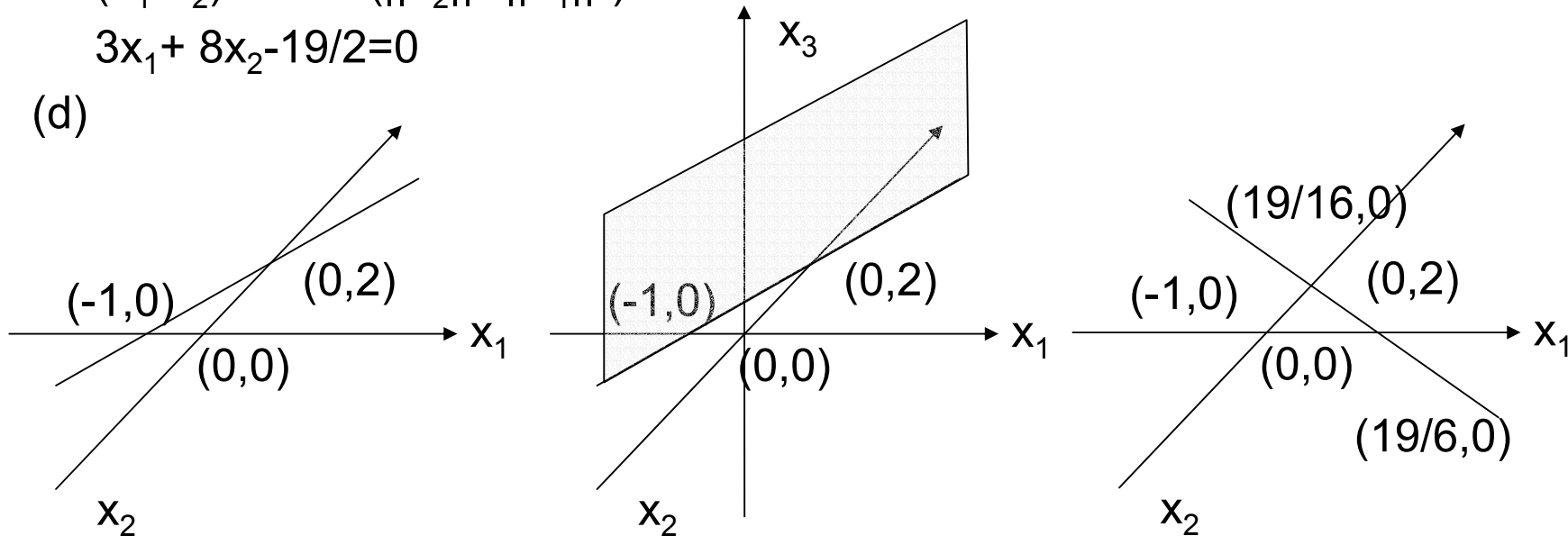
(c) $\mathbf{x}_1 = [2, 5]$, $\mathbf{x}_2 = [-1, -3]$

=> The decision surface for minimum distance classifier

$$(\mathbf{x}_1 - \mathbf{x}_2)^t \mathbf{x} + \frac{1}{2} (\|\mathbf{x}_2\|^2 - \|\mathbf{x}_1\|^2) = 0$$

$$3x_1 + 8x_2 - 19/2 = 0$$

(d)

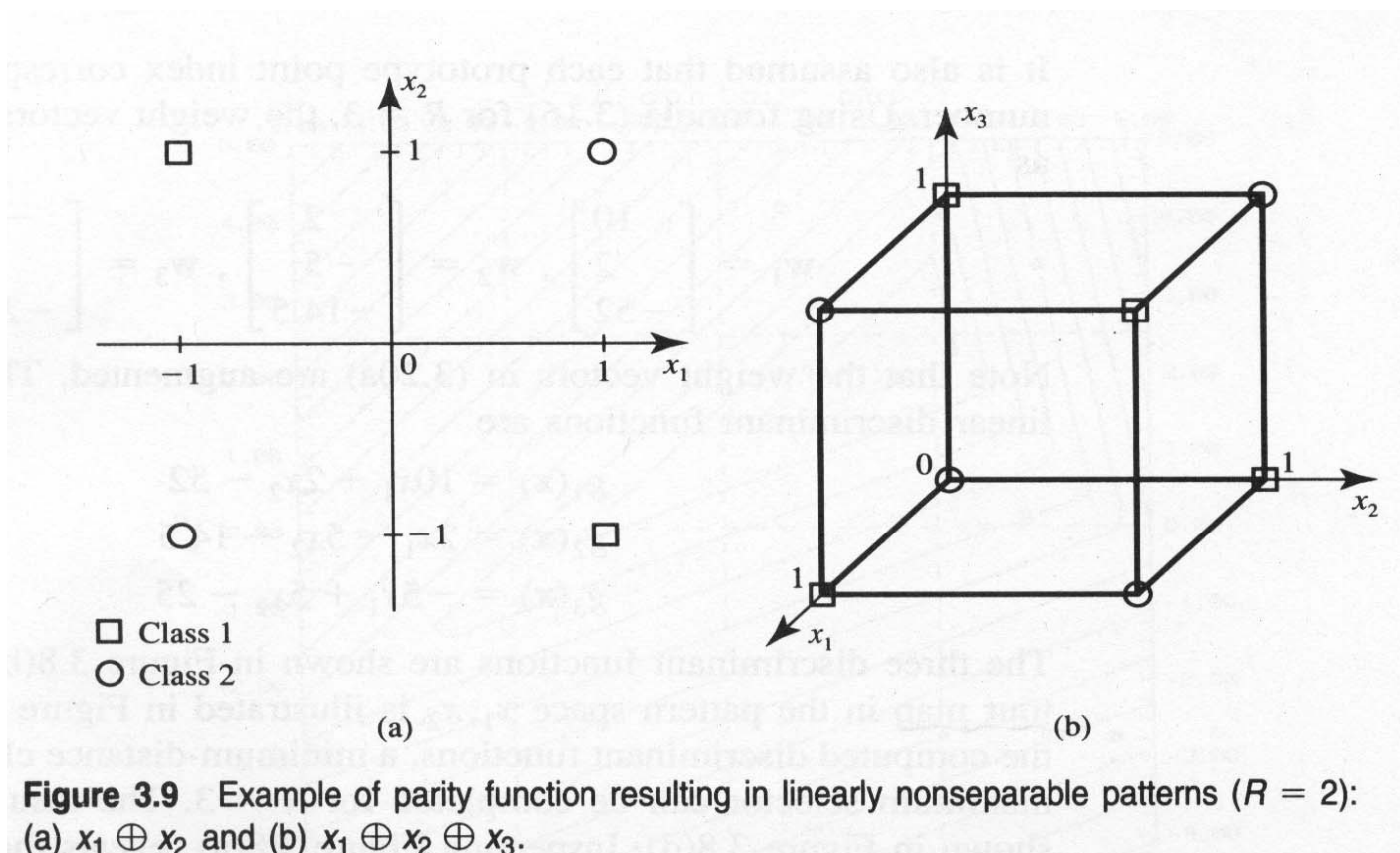


Linear Machine and Minimum Distance Classification (cont.)

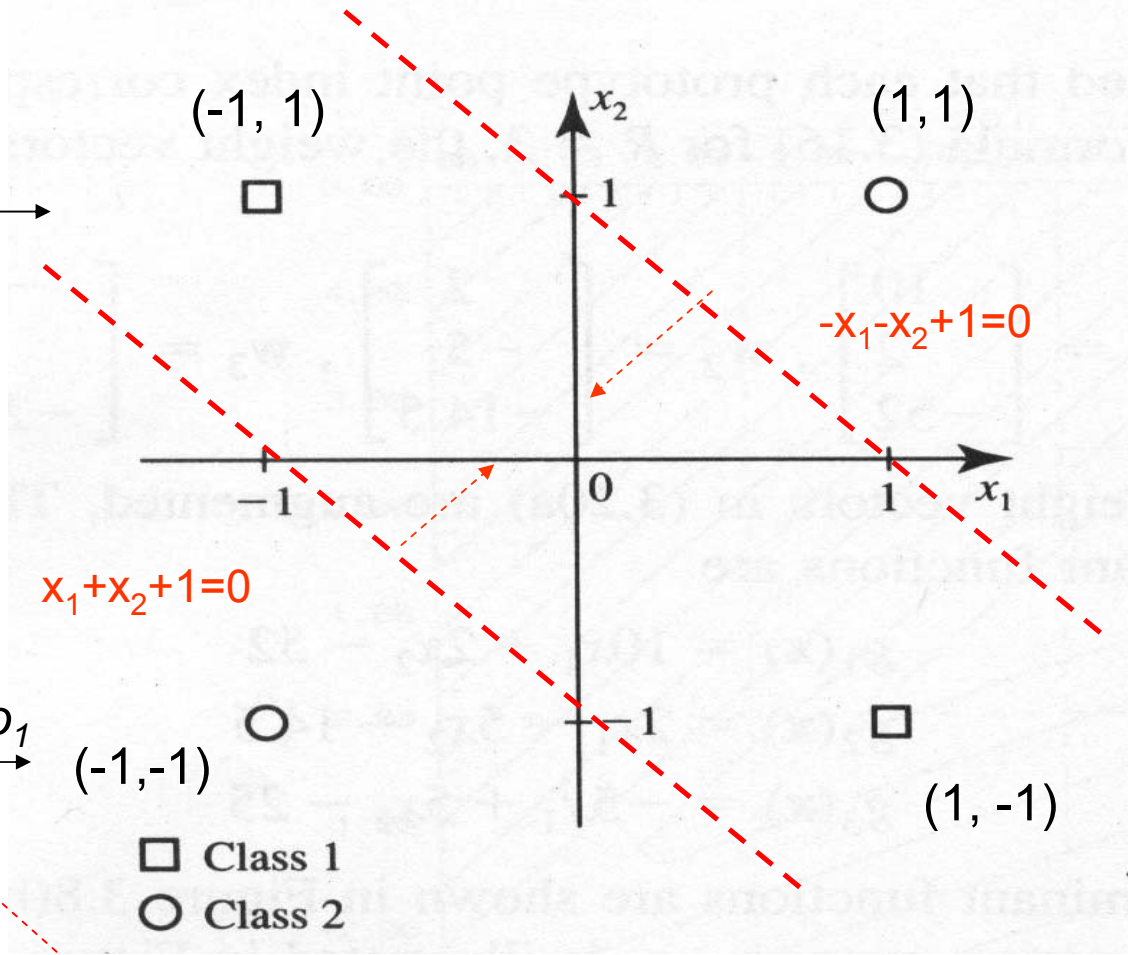
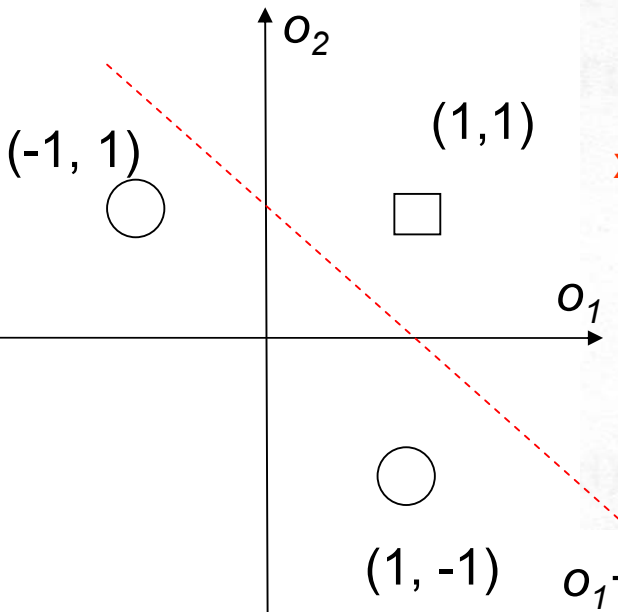
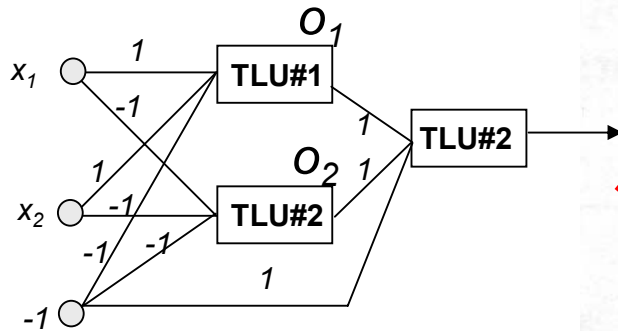
- Examples 3.1 and 3.2 have shown that the coefficients (weights) of the linear discriminant functions can be determined if the a priori information about the sets of patterns and their class membership is known

Linear Machine and Minimum Distance Classification (cont.)

- The example of linearly non-separable patterns



Linear Machine and Minimum Distance Classification (cont.)



Discrete Perceptron Training Algorithm

- Geometrical Representations

- Examine the neural network classifiers that derive/training their weights based on the **error-correction scheme**

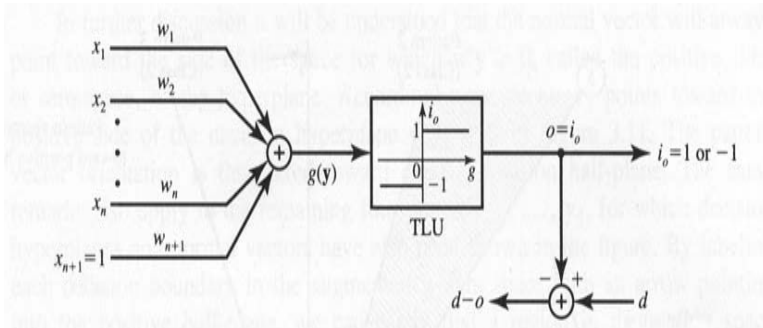


Figure 3.10 Linear dichotomizer using hard-limiting threshold element, or the TLU-based perceptron.

Class 1: $w^t y > 0$

Class 2: $w^t y < 0$

$g(y) = w^t y$

Augmented input pattern

Vector Representations in the **Weight Space**

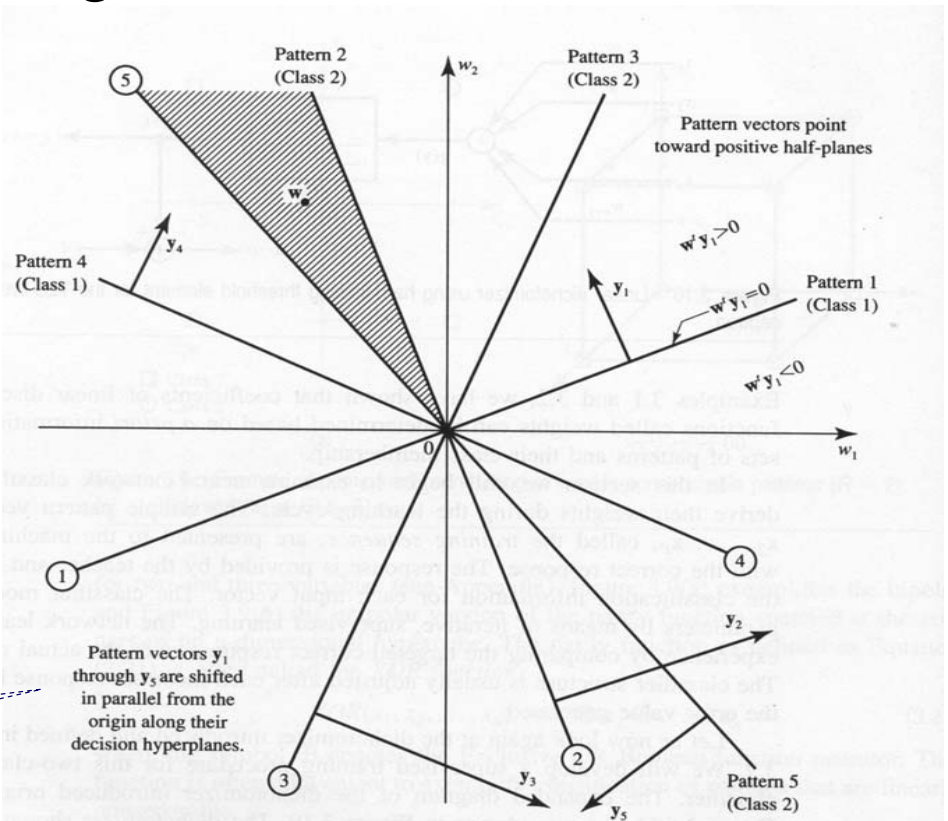
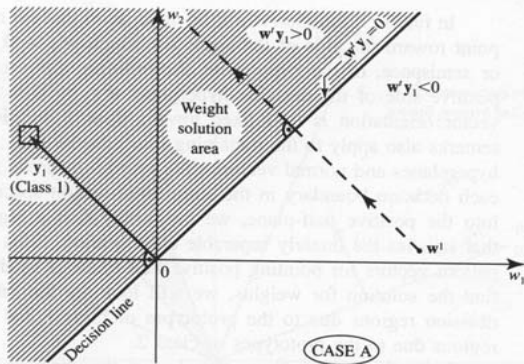


Figure 3.11 Decision hyperplane in augmented weight space for a five pattern set from two classes.

Discrete Perceptron Training Algorithm - Geometrical Representations (cont.)

- Devise an analytic approach based on the geometrical representations
 - E.g. the decision surface for the training pattern y_1

y_1 in Class 1



$$\nabla_w (w^t y_1) = y_1$$

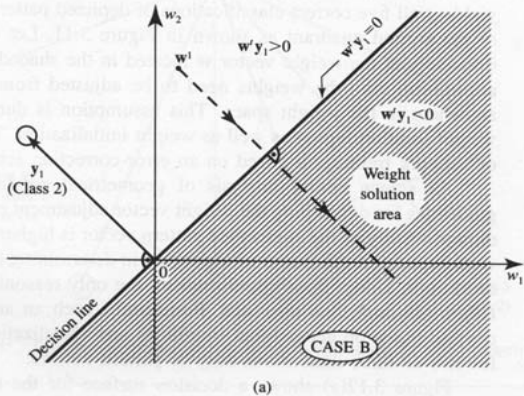
Gradient
(the direction of steep increase)

If y_1 in Class 1:

$$w' = w^1 + cy_1$$

Weight Space

y_1 in Class 2



If y_1 in Class 2:

$$w' = w^1 - cy_1$$

c controls the size of adjustment

$c (>0)$ is the correction increment (is two times of the learning constant introduced before)

Figure 3.12a Weight adjustments of learning dichotomizer: (a) steepest descent

Discrete Perceptron Training Algorithm - Geometrical Representations (count.)

Weight adjustments of three augmented training pattern \mathbf{y}_1 , \mathbf{y}_2 , \mathbf{y}_3 , shown in the weight space

$$\mathbf{y}_1 \in C_1$$

$$\mathbf{y}_2 \in C_1$$

$$\mathbf{y}_3 \in C_2$$

- Weights in the shaded region are the solutions
- The three lines labeled are fixed during training

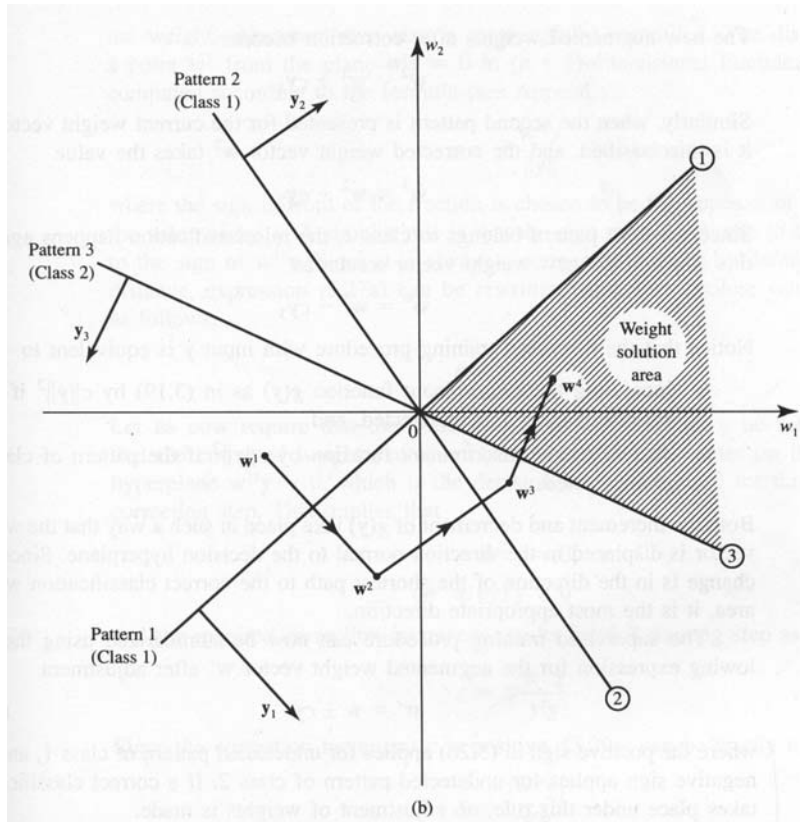
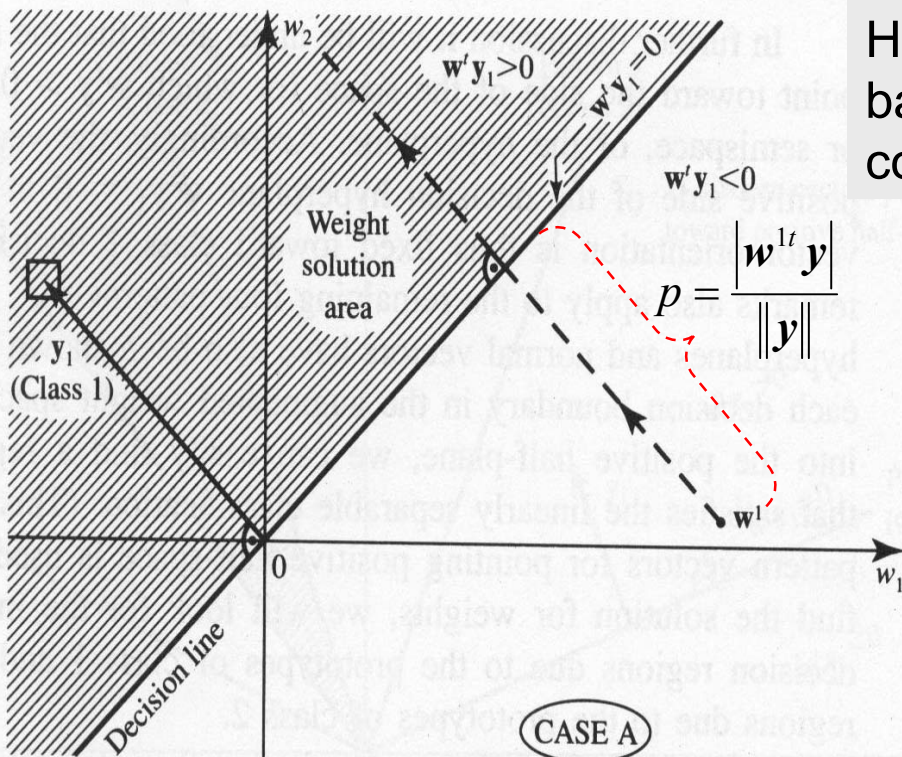


Figure 3.12b Weight adjustments of learning dichotomizer (continued): (b) example.

Discrete Perceptron Training Algorithm

- Geometrical Representations (count.)

- More about the correction increment c
 - If it is not merely a constant, but related to the current training pattern



How to select the correction increment based on the dislocates of w^1 and the corrected weight vector w

$$\begin{aligned}
 (w^1 \pm cy)^t y &= 0 \\
 c &= \mp \frac{w^{1t} y}{y^t y} = \frac{|w^{1t} y|}{\|y\|^2}, \text{ because } c > 0 \\
 \Rightarrow cy &= \frac{|w^{1t} y|}{\|y\|^2} y
 \end{aligned}$$

Discrete Perceptron Training Algorithm

- Geometrical Representations (count.)

- For **fixed correction rule** with $c = \text{constant}$, the correction of weights is always the same fixed portion of the current training vector
 - The weight can be initialized at any value

$$\mathbf{w}' = \mathbf{w} \pm c\mathbf{y} \quad \text{or} \quad \begin{aligned} \mathbf{w}' &= \mathbf{w} + \Delta \mathbf{w} \\ \Delta \mathbf{w} &= c[d - \text{sgn}(\mathbf{w}^t \mathbf{y})]\mathbf{y} \end{aligned}$$

- For **dynamic correction rule** with c dependent on the distance from the weight (i.e. the weight vector) to the decision surface in the weight space

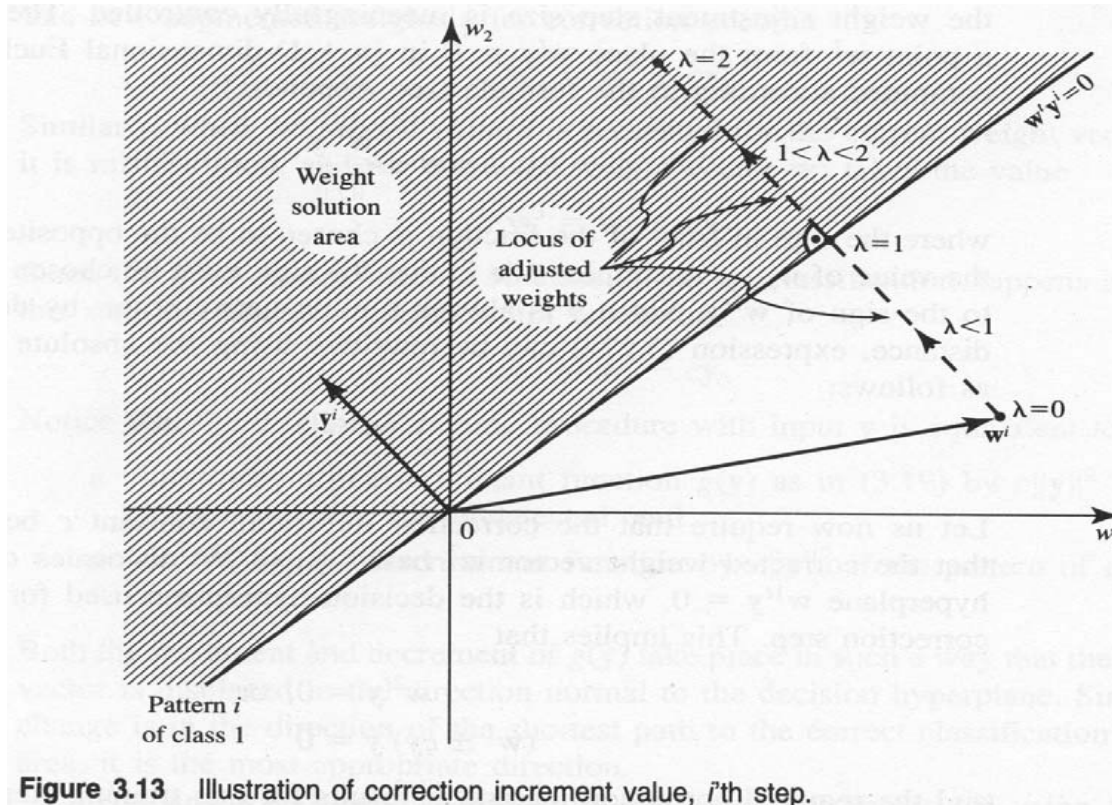
- The initial weight should be different from $\mathbf{0}$

$$\Rightarrow c\mathbf{y} = \frac{|\mathbf{w}^{1t} \mathbf{y}|}{\|\mathbf{y}\|^2} \mathbf{y}$$

Discrete Perceptron Training Algorithm

- Geometrical Representations (count.)

- Dynamic correction rule with c dependent on the distance from the weight



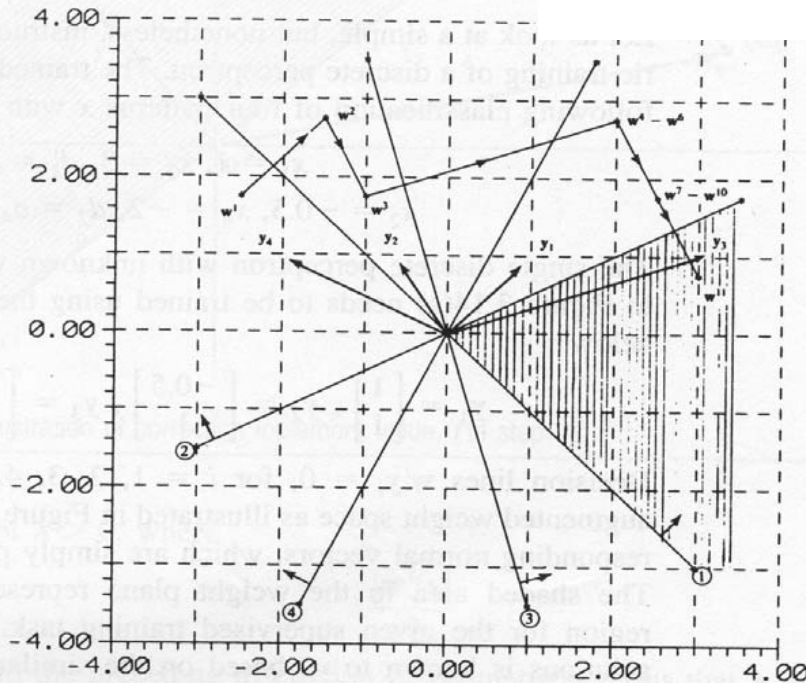
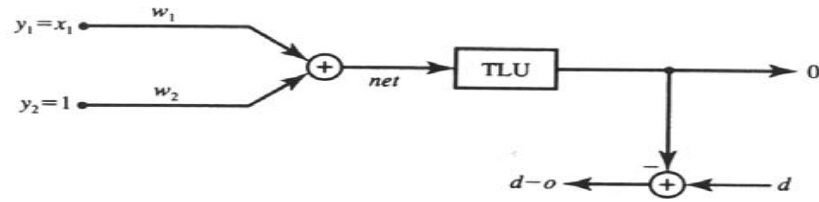
$$c = \lambda \frac{|w^{1t} y|}{\|y\|^2}$$

$$c y = \lambda \frac{|w^{1t} y|}{\|y\|} \frac{y}{\|y\|}$$

Figure 3.13 Illustration of correction increment value, i 'th step.

Discrete Perceptron Training Algorithm - Geometrical Representations (count.)

- Example 3.3



(a)

$$y_1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad y_2 = \begin{bmatrix} -0.5 \\ 1 \end{bmatrix} \quad \begin{array}{l} y_1 \in C_1 \\ y_2 \in C_2 \end{array}$$

$$y_3 = \begin{bmatrix} 3 \\ 1 \end{bmatrix} \quad y_4 = \begin{bmatrix} 2 \\ -1 \end{bmatrix} \quad \begin{array}{l} y_3 \in C_1 \\ y_4 \in C_2 \end{array}$$

$$\Delta w^k = \frac{c}{2} \left[d_k - \text{sgn}(\mathbf{w}^{kt} \mathbf{y}_j) \right] \mathbf{y}_j$$

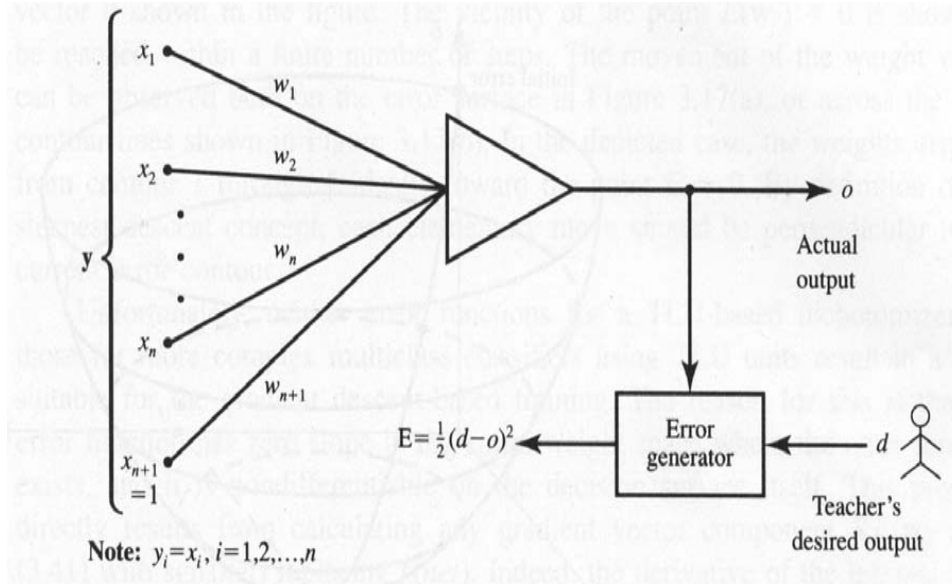
What if $\mathbf{w}^{kt} \mathbf{y}_j = 0$?
 -> interpreted as a mistake
 and followed by a correlation

(b)

Figure 3.14a,b Discrete perceptron classifier training in Example 3.3: (a) network diagram, (b) fixed correction rule training.

Continuous Perceptron Training Algorithm

- Replace the TLU (Threshold Logic Unit) with the sigmoid activation function for two reasons:
 - Gain finer control over the training procedure
 - Facilitate the differential characteristics to enable computation of the error gradient



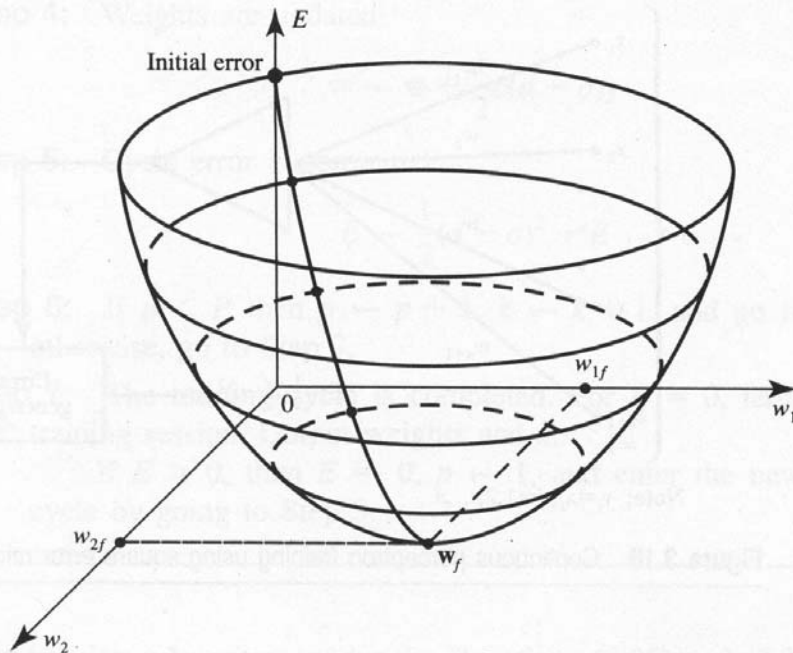
$$\hat{w} = w - \eta \nabla E(w)$$

learning constant error gradient

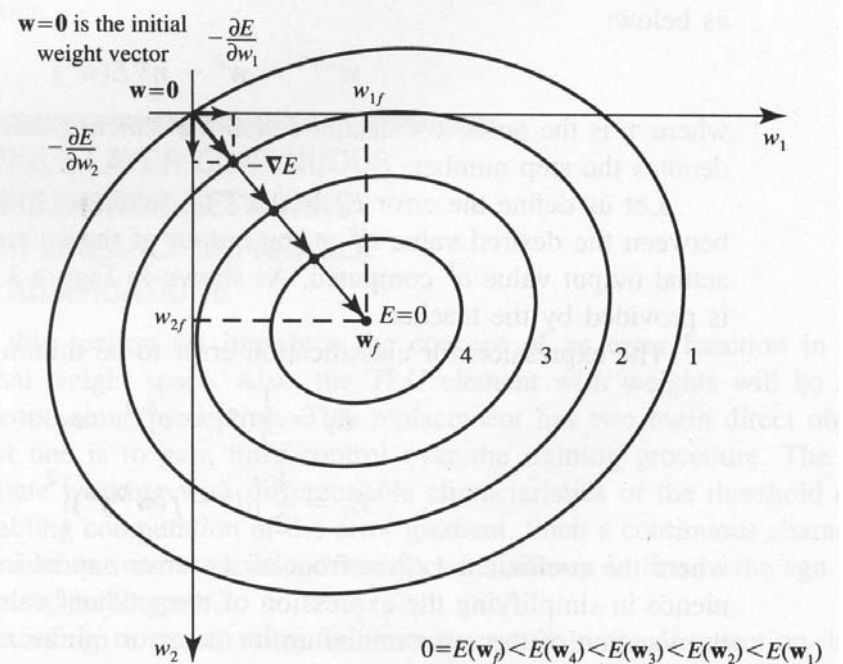
Figure 3.16 Continuous perceptron training using square error minimization.

Continuous Perceptron Training Algorithm (cont.)

- The new weights are obtained by moving in the direction of the negative gradient along the multidimensional error surface



(a)



(b)

Continuous Perceptron Training Algorithm (cont.)

- Define the error as the squared difference between the desired output and the actual output

$$E = \frac{1}{2}(d - o)^2$$

$$\text{or } E = \frac{1}{2}[d - f(\mathbf{w}^t \mathbf{y})]^2 = \frac{1}{2}[d - f(\text{net})]^2$$

$$\nabla E(\mathbf{w}) = \frac{1}{2} \nabla ([d - f(\text{net})]^2)$$

$$\nabla E(\mathbf{w}) = \begin{bmatrix} \frac{\partial E}{\partial w_1} \\ \frac{\partial E}{\partial w_2} \\ \vdots \\ \frac{\partial E}{\partial w_{n+1}} \end{bmatrix} = -(d - o) f'(\text{net}) \begin{bmatrix} \frac{\partial (\text{net})}{\partial w_1} \\ \frac{\partial (\text{net})}{\partial w_2} \\ \vdots \\ \frac{\partial (\text{net})}{\partial w_{n+1}} \end{bmatrix} = -(d - o) f'(\text{net}) \mathbf{y}$$

Continuous Perceptron Training Algorithm (cont.)

- Bipolar Continuous Activation Function

$$f(\text{net}) = \frac{2}{1 + \exp(-\lambda \cdot \text{net})} - 1 \quad f'(\text{net}) = \lambda \cdot \frac{2 \exp(-\lambda \cdot \text{net})}{[1 + \exp(-\lambda \cdot \text{net})]^2} = \lambda \cdot \{1 - [f(\text{net})]^2\} = \lambda(1 - o^2)$$

$$\hat{\mathbf{w}} = \mathbf{w} + \frac{1}{2} \eta \cdot \lambda (d - o)(1 - o^2) \mathbf{y}$$

- Unipolar Continuous Activation Function

$$f(\text{net}) = \frac{1}{1 + \exp(-\lambda \cdot \text{net})} \quad f'(\text{net}) = \frac{\lambda \cdot \exp(-\lambda \cdot \text{net})}{[1 + \exp(-\lambda \cdot \text{net})]^2} = \lambda \cdot f(\text{net})[1 - f(\text{net})] = \lambda \cdot o(1 - o)$$

$$\hat{\mathbf{w}} = \mathbf{w} + \eta \cdot \lambda \cdot (d - o) o (1 - o) \mathbf{y}$$

Continuous Perceptron Training Algorithm (cont.)

- Example 3.3 $f(net) = \frac{2}{1 + \exp(-net)} - 1$

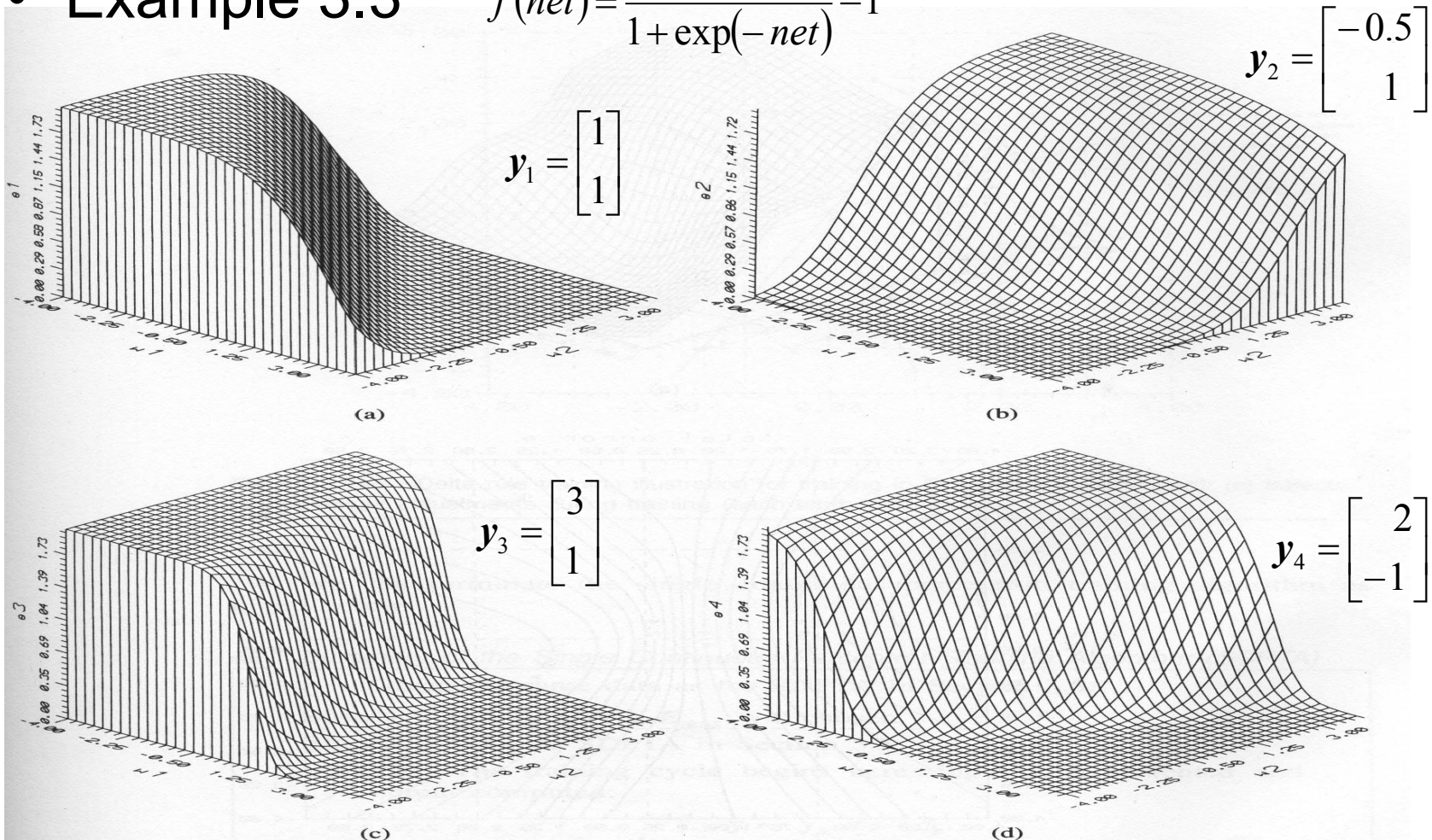


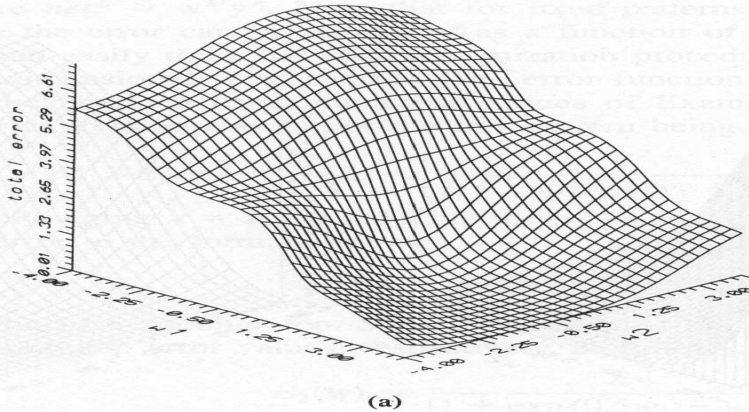
Figure 3.18 Error functions for individual patterns, Example 3.4: (a) first pattern, E_1 , (b) second pattern, E_2 , (c) third pattern, and E_3 , and (d) fourth pattern, E_4 .

Continuous Perceptron Training Algorithm (cont.)

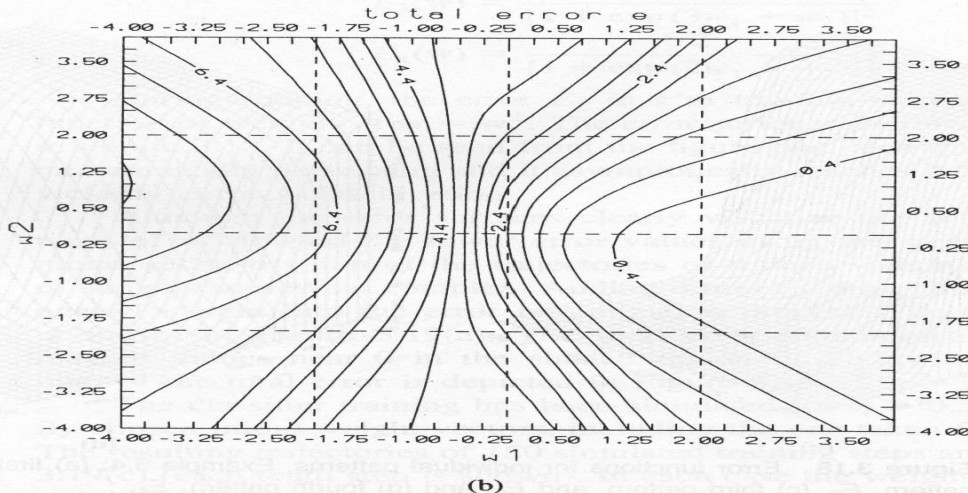
- Example 3.3

Total error surface

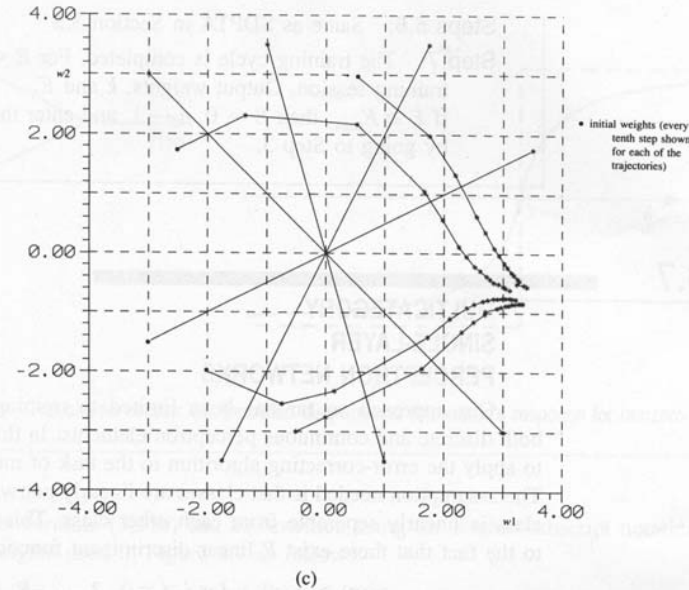
Trajectories started from four arbitrary initial weights



(a)



(b)



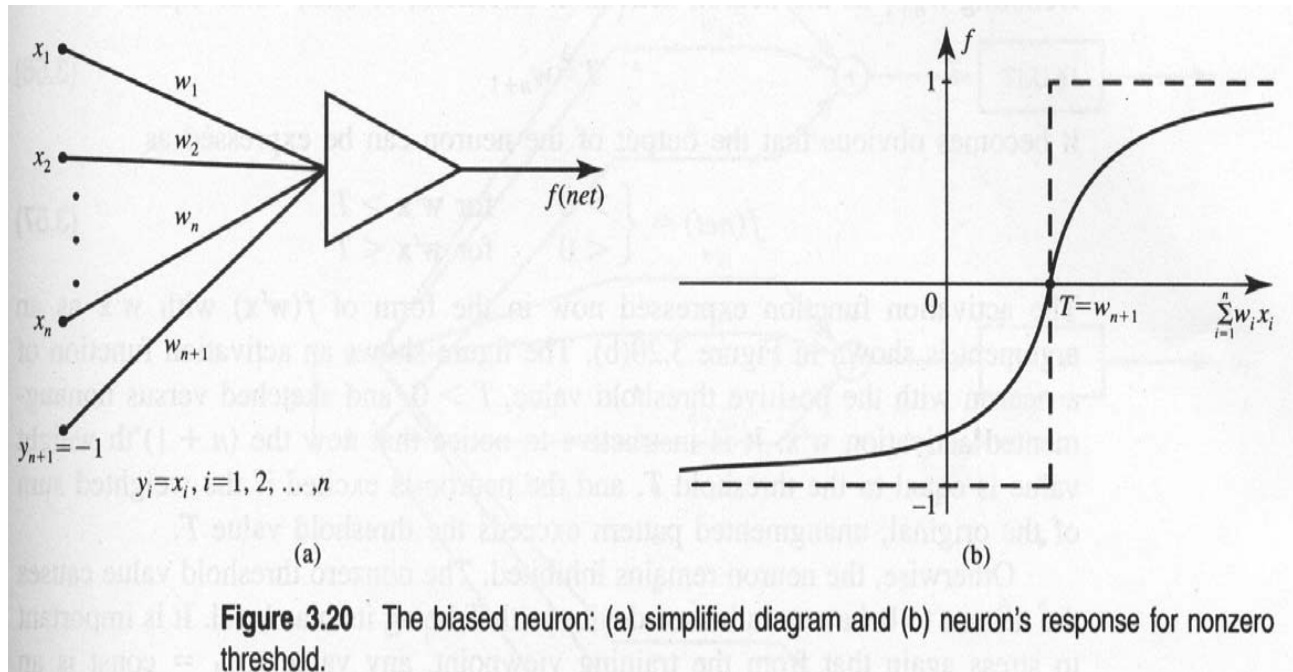
(c)

Figure 3.19c Delta rule training illustration for training in Example 3.4 (continued): (c) trajectories of weight adjustments during training (each tenth step shown).

Figure 3.19a,b Delta rule training illustration for training in Example 3.4: (a) total error surface, (b) total error contour map.

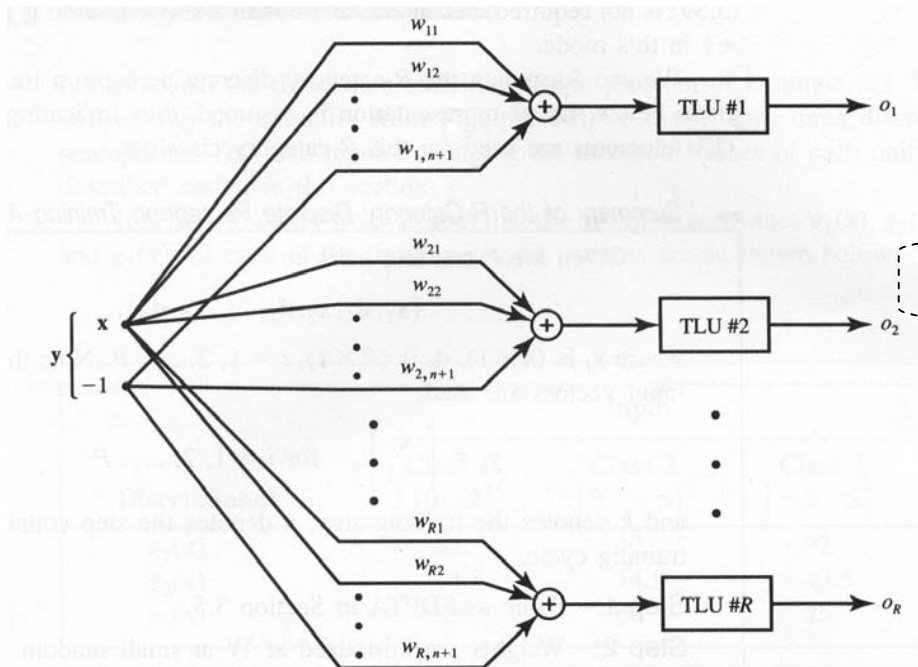
Continuous Perceptron Training Algorithm (cont.)

- Treat the last fixed component of input pattern vector as the neuron activation threshold



Continuous Perceptron Training Algorithm (cont.)

- R -category linear classifier using R discrete bipolar perceptrons
 - Goal: The i -th TLU response of +1 is indicative of class i and all other TLU respond with -1



$$\hat{\mathbf{w}}_i = \mathbf{w}_i + \frac{1}{2} c \cdot (d_i - o_i) \mathbf{y}$$

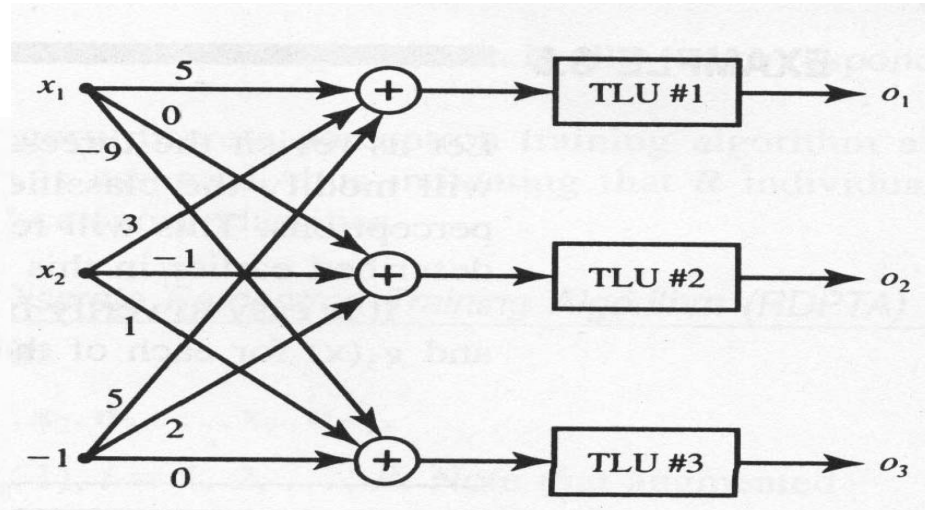
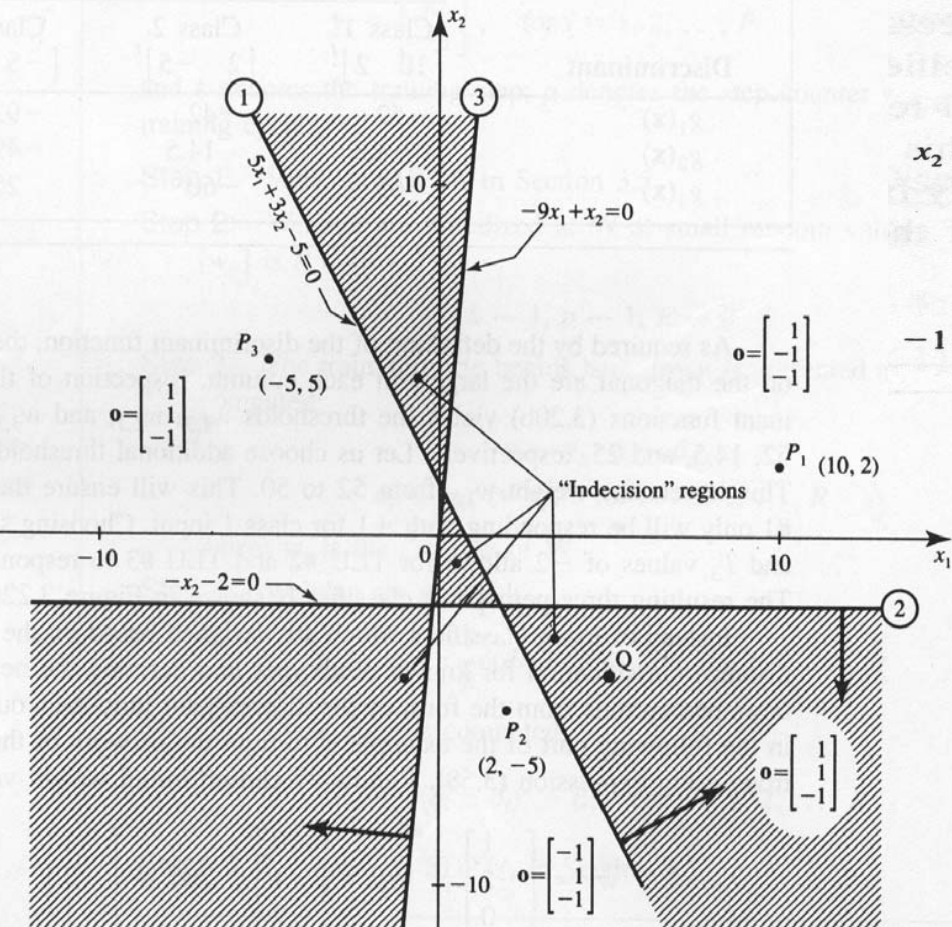
$$d_i = 1, d_j = -1, \text{ for } j = 1, 2, \dots, R, j \neq i$$

For “local representation”

Figure 3.21 R -category linear classifier using R discrete perceptrons.

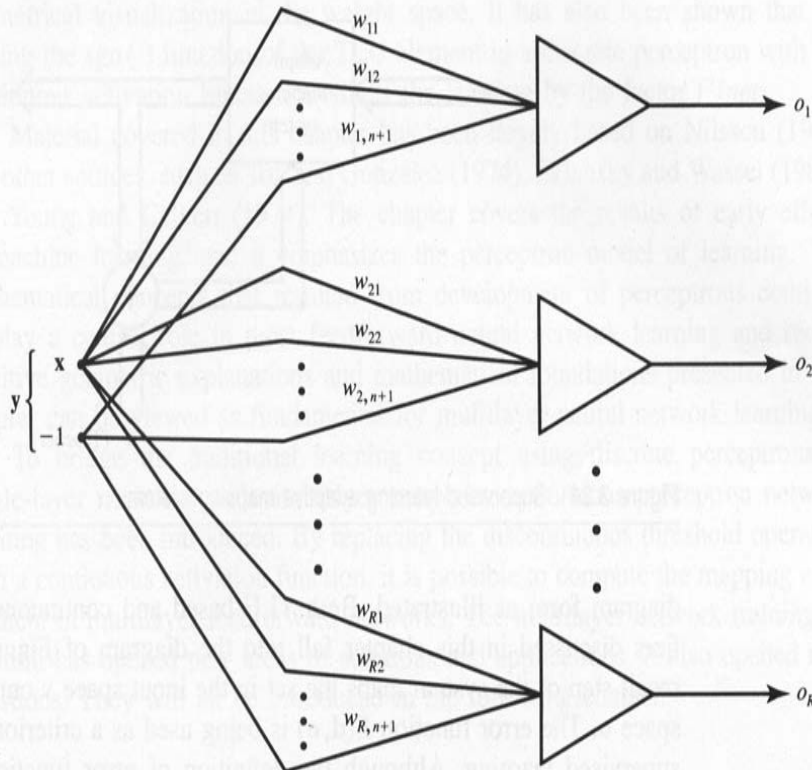
Continuous Perceptron Training Algorithm (cont.)

- Example 3.5



Continuous Perceptron Training Algorithm (cont.)

- R -category linear classifier using R continuous bipolar perceptrons



$$\hat{\mathbf{w}}_i = \mathbf{w}_i + \frac{1}{2} \eta \cdot \lambda (d_i - o_i) (1 - o_i^2) \mathbf{y}$$

for $i = 1, 2, \dots, R$

$$d_i = 1, d_j = -1, \text{ for } j = 1, 2, \dots, R, j \neq i$$

Figure 3.23 R -category linear classifier using continuous perceptrons.

Continuous Perceptron Training Algorithm (cont.)

- Error function dependent on the difference vector $\mathbf{d}-\mathbf{o}$

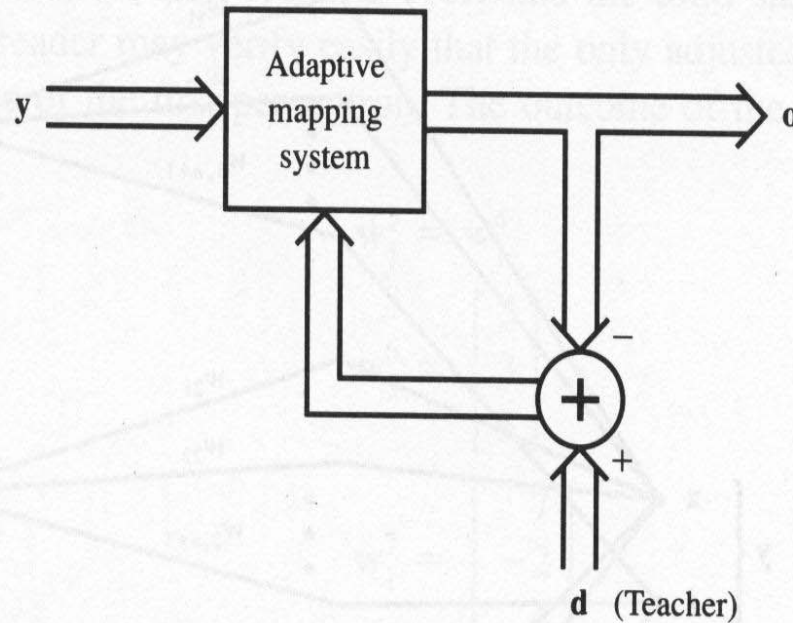


Figure 3.24 Supervised learning adaptive mapping system.

Bayes' Classifier vs. Perceptron

- Perceptron operates on the promise that the patterns to be classified are *linear separable* (otherwise the training algorithm will oscillate), while Bayes' classifier assumes the (Gaussian) distribution of two classes certainly do overlap each other
- The perceptron is nonparametric while the Bayes' classifier is parametric (its derivation is contingent on the assumption of the underlying distributions)
- The perceptron is simple and adaptive, and needs small storage, while the Bayes' classifier could be made adaptive but at the expense of increased storage and more complex computations

Homework

- P3.5, P3.7, P3.9, P3.22