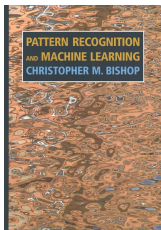


# Linear Models for Classification

Bishop PRML Ch. 4 – 4.1

Asahi Tsuruo

June 3, 2020



## 1 4.1 Discriminant Functions

- 4.1.1 Two classes
- 4.1.2 Multiple classes
- 4.1.3 Least squares for classification
- 4.1.4 Fisher's linear discriminant
- 4.1.5 Relation to least squares
- 4.1.6 Fisher's discriminant for multiple classes
- 4.1.7 The perceptron algorithm

# Discriminant functions

Example in Hand-written Digit Recognition:

input vector

target vector

$$\mathbf{x}_i = \begin{array}{|c|} \hline \text{7} \\ \hline \end{array} \quad \mathbf{t}_i = \{0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0\} \quad (4.1^*)$$

- Each input vector  $\mathbf{x}_i$  is classified into one of the  $K$  discrete classes
  - Denote classes  $C_k$
- Represent input image as a vector  $\mathbf{x}_i \in \mathbb{R}^{784}$
- Target vector  $\mathbf{t}_i \in \mathbb{R}^{10}$ 
  - 1-of- $K$
- Given a training data set  $\{(\mathbf{x}_n, \mathbf{t}_n)\}_{n=1}^{n=N}$ , learning problem is to construct a "good" function  $y(\mathbf{x})$  from these:
  - $y : \mathbb{R}^{784} \rightarrow \mathbb{R}^{10}$

# Generalized linear models

- Similar to previous chapter on linear models for regression, we'll use a "linear" model for classification:

$$y(\mathbf{x}) = f(\mathbf{w}^T \mathbf{x} + w_0) \quad (4.3)$$

- This (4.3) is called a generalized linear model
- $f(a)$  is an activate function (a nonlinear function)
  - e.g.

$$f(a) = \begin{cases} +1, & a \geq 0 \\ -1, & a < 0 \end{cases} \quad (4.53)$$

- Decision boundary between classes will be linear function of  $\mathbf{x}$
- Can also apply non-linearity to  $\mathbf{x}$ , as in  $\phi_i(\mathbf{x})$  for regression

# Outline

## 1 4.1 Discriminant Functions

- 4.1.1 Two classes
- 4.1.2 Multiple classes
- 4.1.3 Least squares for classification
- 4.1.4 Fisher's linear discriminant
- 4.1.5 Relation to least squares
- 4.1.6 Fisher's discriminant for multiple classes
- 4.1.7 The perceptron algorithm

# Discriminant Functions in binary classification

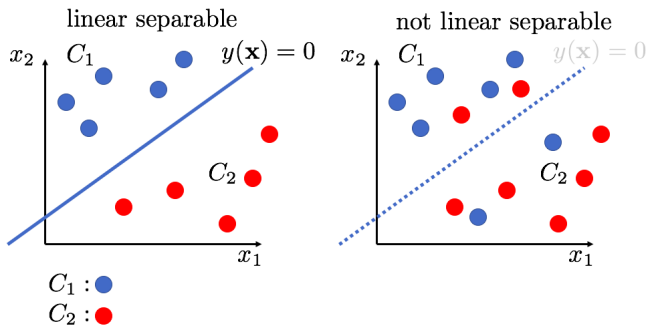
- A simple linear discriminant function:

$$\begin{aligned}y(\mathbf{x}) &= \sum_{j=1}^D w_j x_j + w_0 \\ &= \mathbf{w}^T \mathbf{x} + w_0\end{aligned}\tag{4.4}$$

$\mathbf{w}$ : weight vector,  $\mathbf{x}$ : input vector,  $w_0$ : bias,  $D$ : input space dimension

- An input vector  $\mathbf{x}$  is assigned to
  - class  $C_1$  if  $y(\mathbf{x}) \geq 0$
  - class  $C_2$  otherwise
- Decision boundary:
  - $y(\mathbf{x}) = 0$
  - $(D - 1)$  dimensional hyperplane within input space  $D$

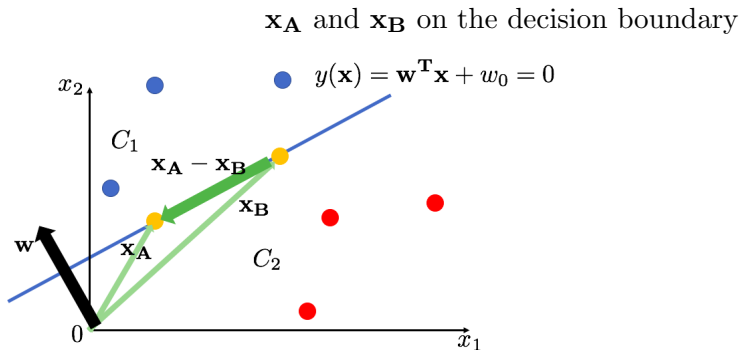
# Example of linear separable and $D - 1$ hyperplane



- Input space (planes):  $2 D$
- Decision boundary (lines):  $(2 - 1) D = 1 D$

If input space is  $3D$  (solids), decision boundary is  $2D$  (planes)

$\mathbf{w}$  determines the orientation of the decision boundary



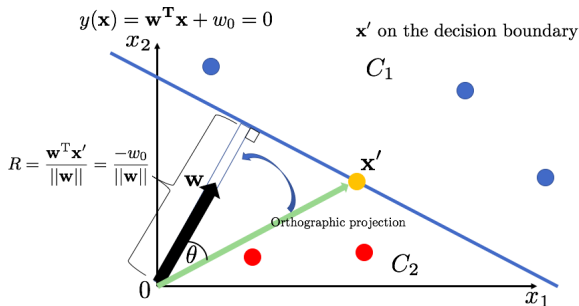
$$y(\mathbf{x}_A) - y(\mathbf{x}_B) = \mathbf{w}^T (\mathbf{x}_A - \mathbf{x}_B) = 0 \quad (\text{A1})$$

- From an inner product property,  $\mathbf{w}$  is orthogonal to every vector on the decision boundary
  - $\mathbf{w}$  determines the orientation of the decision boundary



# The distance $R$ from the origin to the decision boundary

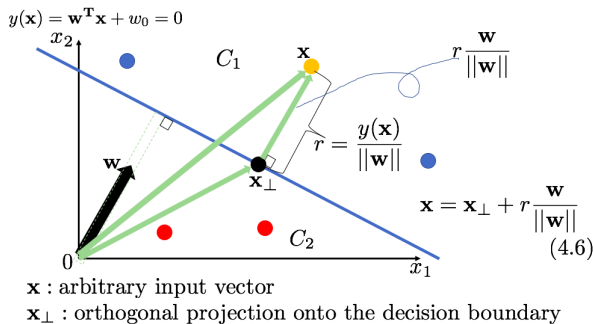
- Bias parameter  $w_0$  determines the location of the decision boundary



$$R = \|\mathbf{e}\| \|\mathbf{x}'\| \cos \theta = \mathbf{e}^T \cdot \mathbf{x}' = \frac{\mathbf{w}^T \mathbf{x}'}{\|\mathbf{w}\|} = \frac{-w_0}{\|\mathbf{w}\|} \quad (4.5^*)$$

- $\mathbf{w}^T / \|\mathbf{w}\| = \mathbf{e}$  is a unit vector,  $\mathbf{w}^T \mathbf{x}' + w_0 = 0 \rightarrow \mathbf{w}^T \mathbf{x}' = -w_0$

The perpendicular distance  $r$  between a  $\mathbf{x}$  and the decision boundary



Multiplying both sides of this result (4.6) by  $\mathbf{w}^T$  and adding  $w_0$ , and making use of  $y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$  and  $y(\mathbf{x}_\perp) = \mathbf{w}^T \mathbf{x}_\perp + w_0 = 0$

$$r = \frac{y(\mathbf{x})}{\|\mathbf{w}\|} \quad (4.7)$$

# Augmented representation for binary classification

■  $x_0 = 1, \hat{\mathbf{w}} = (w_0, \mathbf{w}), \hat{\mathbf{x}} = (x_0, \mathbf{x})$

$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 \quad (4.4)$$

$$= \mathbf{w}^T \mathbf{x} + w_0 * 1$$

$$= \mathbf{w}^T \mathbf{x} + w_0 * x_0$$

$$= \tilde{\mathbf{w}}^T \tilde{\mathbf{x}} \quad (4.8)$$

- The decision boundaries are  $D$ -dimensional hyperplanes passing through the origin of the  $D + 1$ -dimensional expanded input space

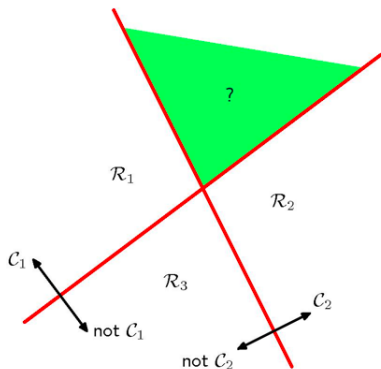
# Outline

## 1 4.1 Discriminant Functions

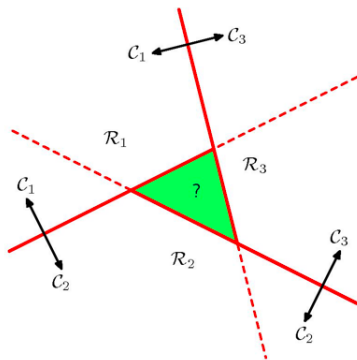
- 4.1.1 Two classes
- 4.1.2 Multiple classes
- 4.1.3 Least squares for classification
- 4.1.4 Fisher's linear discriminant
- 4.1.5 Relation to least squares
- 4.1.6 Fisher's discriminant for multiple classes
- 4.1.7 The perceptron algorithm

# Discriminant Functions in Multiple Classes

One-versus-the-rest classifier



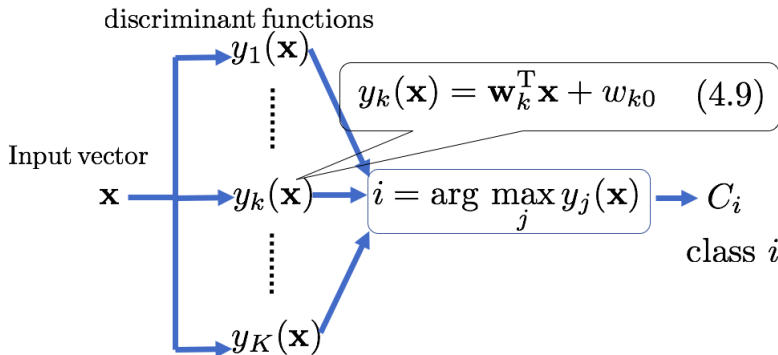
One-versus-one classifier



- A linear discriminant function between two classes separates with a hyperplane
- How can we use this for multiple classes?
  - One-versus-the-rest classifier: build  $K - 1$  classifiers
  - One-versus-one classifier: build  $K(K - 1)/2$

# Solution for ambiguous regions in Multiple Classes

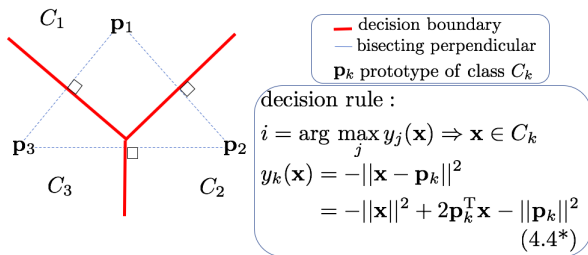
- Select the discriminant function that have a maximum output



- The decision boundary between  $C_k$  and  $C_j$ 
  - Given by  $y_k(\mathbf{x}) = y_j(\mathbf{x})$
  - Corresponds to a  $D - 1$  dimensional hyperplane

$$(\mathbf{w}_k - \mathbf{w}_j)^T \mathbf{x} + (w_{k0} - w_{j0}) = 0 \quad (4.10)$$

Nearest Neighbour (NN) classifier is a specific case of linear discriminants in selection of max-value



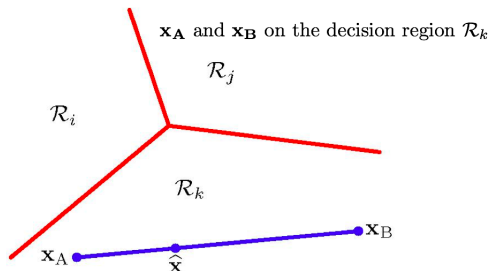
From NN classifier (4.4\*), we can obtain the classifier  $\mathbf{p}_k^T \mathbf{x} - \frac{1}{2}\|\mathbf{p}_k\|^2$ , thus,

$$y(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0} \quad (4.4)$$

where

$$\mathbf{w}_k = \mathbf{p}_k, \quad w_0 = -\frac{1}{2}\|\mathbf{p}_k\|^2 \quad (\text{B1})$$

# The decision regions are connected and convex



- $\hat{\mathbf{x}}$ : any point that lies on the line connecting  $\mathbf{x}_A$  and  $\mathbf{x}_B$
- $0 \leq \lambda \leq 1$

$$\hat{\mathbf{x}} = \lambda \mathbf{x}_A + (1 - \lambda) \mathbf{x}_B \quad (4.11)$$

$$y_k(\hat{\mathbf{x}}) = \lambda y_k(\mathbf{x}_A) + (1 - \lambda) y_k(\mathbf{x}_B) \quad (4.12)$$

$$\Rightarrow y_k(\hat{\mathbf{x}}) > y_j(\hat{\mathbf{x}}), \forall j \neq k$$



# Outline

## 1 4.1 Discriminant Functions

- 4.1.1 Two classes
- 4.1.2 Multiple classes
- 4.1.3 Least squares for classification
- 4.1.4 Fisher's linear discriminant
- 4.1.5 Relation to least squares
- 4.1.6 Fisher's discriminant for multiple classes
- 4.1.7 The perceptron algorithm

# Augmented representation for multi-class classification

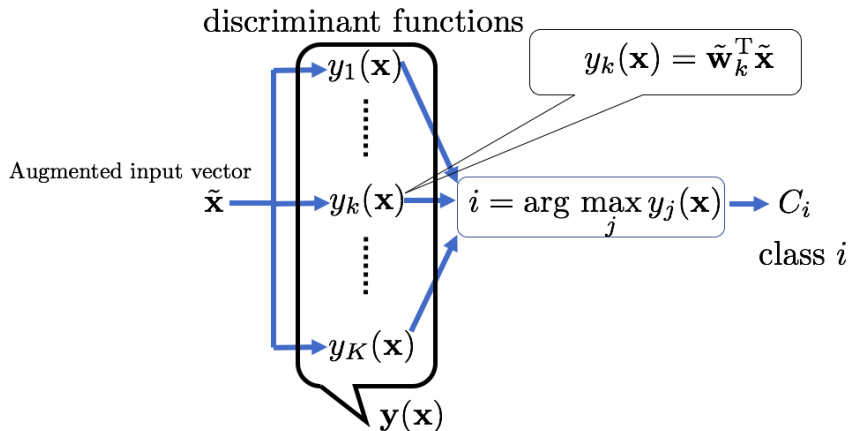
- The discriminant functions for each class  $C_k$

$$y_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0} \quad (k = 1, \dots, K) \quad (4.13)$$

$\Downarrow$

$$\mathbf{y}(\mathbf{x}) = \widetilde{\mathbf{W}}^T \tilde{\mathbf{x}} \quad (4.14)$$

- $\widetilde{\mathbf{W}}$  is a matrix whose  $k^{\text{th}}$  column consists of the  $D + 1$  dimensional vector  $\tilde{\mathbf{w}}_k = (w_{k0}, \mathbf{w}_k^T)^T$
- $\tilde{\mathbf{x}}$  is the augmented input vector  $(1, \mathbf{x}^T)^T$



$$\mathbf{y}(\mathbf{x}) = \tilde{\mathbf{W}}^T \tilde{\mathbf{x}} = \begin{pmatrix} w_{10} & \cdots & w_{1d} & \cdots & w_{1D} \\ \vdots & \ddots & & & \vdots \\ w_{k0} & & \ddots & & w_{kD} \\ \vdots & & & \ddots & \vdots \\ w_{K0} & \cdots & w_{Kd} & \cdots & w_{KD} \end{pmatrix} \begin{pmatrix} x_0 \\ \vdots \\ x_d \\ \vdots \\ x_D \end{pmatrix} = \begin{pmatrix} \tilde{\mathbf{w}}_1^T \tilde{\mathbf{x}} \\ \vdots \\ \tilde{\mathbf{w}}_k^T \tilde{\mathbf{x}} \\ \vdots \\ \tilde{\mathbf{w}}_K^T \tilde{\mathbf{x}} \end{pmatrix} = \begin{pmatrix} y_1(\mathbf{x}) \\ \vdots \\ y_k(\mathbf{x}) \\ \vdots \\ y_K(\mathbf{x}) \end{pmatrix}$$

# How do we learn the decision boundaries $(w_0, \mathbf{w})$ ?

- One approach is to use least squares, similar to regression

$$\begin{aligned} E_D(\widetilde{\mathbf{W}}) &= \frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K \{y_k(\mathbf{x}_n) - t_{nk}\}^2 \\ &= \frac{1}{2} \text{Tr} \left\{ (\widetilde{\mathbf{X}}\widetilde{\mathbf{W}} - \mathbf{T})^T (\widetilde{\mathbf{X}}\widetilde{\mathbf{W}} - \mathbf{T}) \right\} \end{aligned} \quad (4.15)$$

- $\widetilde{\mathbf{X}}$  is a matrix whose  $n^{th}$  row is  $\tilde{\mathbf{x}}_n^T$
- $\mathbf{T}$  is a matrix whose  $n^{th}$  row is the vector  $\mathbf{t}_n^T$

# Let's get used to a matrix representation

$$\begin{aligned}
 E_D(\tilde{\mathbf{W}}) &= \frac{1}{2} \text{Tr} \left\{ (\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{T})^T (\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{T}) \right\} \\
 &= \frac{1}{2} \text{Tr} \left[ \left\{ \begin{pmatrix} x_{10} & \cdots & x_{1D} \\ \vdots & \ddots & \vdots \\ x_{N0} & \cdots & x_{ND} \end{pmatrix} \begin{pmatrix} w_{10} & \cdots & w_{K0} \\ \vdots & \ddots & \vdots \\ w_{1D} & \cdots & w_{KD} \end{pmatrix} - \begin{pmatrix} t_{11} & \cdots & t_{1K} \\ \vdots & \ddots & \vdots \\ t_{N1} & \cdots & t_{NK} \end{pmatrix} \right\}^T \right. \\
 &\quad \left. \left\{ \begin{pmatrix} x_{10} & \cdots & x_{1D} \\ \vdots & \ddots & \vdots \\ x_{N0} & \cdots & x_{ND} \end{pmatrix} \begin{pmatrix} w_{10} & \cdots & w_{K0} \\ \vdots & \ddots & \vdots \\ w_{1D} & \cdots & w_{KD} \end{pmatrix} - \begin{pmatrix} t_{11} & \cdots & t_{1K} \\ \vdots & \ddots & \vdots \\ t_{N1} & \cdots & t_{NK} \end{pmatrix} \right\} \right] \\
 &= \frac{1}{2} \text{Tr} \left\{ \begin{pmatrix} \tilde{\mathbf{w}}_1^T \tilde{\mathbf{x}}_1 - t_{11} & \cdots & \tilde{\mathbf{w}}_K^T \tilde{\mathbf{x}}_1 - t_{1K} \\ \vdots & \ddots & \vdots \\ \tilde{\mathbf{w}}_1^T \tilde{\mathbf{x}}_N - t_{N1} & \cdots & \tilde{\mathbf{w}}_K^T \tilde{\mathbf{x}}_N - t_{NK} \end{pmatrix}^T \begin{pmatrix} \tilde{\mathbf{w}}_1^T \tilde{\mathbf{x}}_1 - t_{11} & \cdots & \tilde{\mathbf{w}}_K^T \tilde{\mathbf{x}}_1 - t_{1K} \\ \vdots & \ddots & \vdots \\ \tilde{\mathbf{w}}_1^T \tilde{\mathbf{x}}_N - t_{N1} & \cdots & \tilde{\mathbf{w}}_K^T \tilde{\mathbf{x}}_N - t_{NK} \end{pmatrix} \right\} \\
 &= \frac{1}{2} \text{Tr} \begin{pmatrix} (\tilde{\mathbf{w}}_1^T \tilde{\mathbf{x}}_1 - t_{11})^2 + \cdots & , \cdots , & (\tilde{\mathbf{w}}_1^T \tilde{\mathbf{x}}_1 - t_{11})(\tilde{\mathbf{w}}_K^T \tilde{\mathbf{x}}_1 - t_{1K}) + \cdots \\ \vdots & \ddots & \vdots \\ (\tilde{\mathbf{w}}_K^T \tilde{\mathbf{x}}_1 - t_{1K})(\tilde{\mathbf{w}}_1^T \tilde{\mathbf{x}}_1 - t_{11}) + \cdots & , \cdots , & (\tilde{\mathbf{w}}_K^T \tilde{\mathbf{x}}_1 - t_{1K})^2 + \cdots \end{pmatrix} \\
 &= \frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K (y_k(\mathbf{x}_n) - t_{nk})^2
 \end{aligned} \tag{4.15}$$

# Solution for $\widetilde{\mathbf{W}}$

- Setting the derivative with respect to  $\widetilde{\mathbf{W}}$  to zero, and rearranging

$$\widetilde{\mathbf{W}} = (\widetilde{\mathbf{X}}^T \widetilde{\mathbf{X}})^{-1} \widetilde{\mathbf{X}}^T \mathbf{T} = \widetilde{\mathbf{X}}^\dagger \mathbf{T} \quad (4.16)$$

- $\widetilde{\mathbf{X}}^\dagger$  is the pseudo-inverse of the matrix  $\widetilde{\mathbf{X}}$

We get the discriminant functions from (4.16)

$$\mathbf{y}(\mathbf{x}) = \widetilde{\mathbf{W}}^T \tilde{\mathbf{x}} = \mathbf{T}^T \left( \widetilde{\mathbf{X}}^\dagger \right)^T \tilde{\mathbf{x}} \quad (4.17)$$

# Solution for $\widetilde{\mathbf{W}}$ : Calculation process

- $\mathbf{t}_k$  is a column component of the matrix  $\mathbf{T}$ 
  - re-writing the error function and differentiating with respect to  $\mathbf{w}_i$

$$E_D(\widetilde{\mathbf{W}}) = \frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K \{y_k(\mathbf{x}_n) - t_{nk}\}^2 = \frac{1}{2} \sum_{k=1}^K \|\widetilde{\mathbf{X}}\widetilde{\mathbf{w}}_k - \mathbf{t}_k\|^2 \quad (4.15^*)$$

$$\frac{\partial E_D}{\partial \mathbf{w}_i} = \widetilde{\mathbf{X}}^T (\widetilde{\mathbf{X}}\widetilde{\mathbf{w}}_i - \mathbf{t}_i) = 0 \quad (\text{C1})$$

$\Downarrow$

$$\widetilde{\mathbf{w}}_i = (\widetilde{\mathbf{X}}^T \widetilde{\mathbf{X}})^{-1} \widetilde{\mathbf{X}}^T \mathbf{t}_i = \widetilde{\mathbf{X}}^\dagger \mathbf{t}_i \quad (\text{C2})$$

$\Downarrow$

$$\widetilde{\mathbf{W}} = (\widetilde{\mathbf{X}}^T \widetilde{\mathbf{X}})^{-1} \widetilde{\mathbf{X}}^T \mathbf{T} = \widetilde{\mathbf{X}}^\dagger \mathbf{T} \quad (4.16)$$

## An interesting property of least-squares solutions with multiple target variables

- If every target vector in the training set satisfies some linear constraint for some constants  $\mathbf{a}$  and  $b$ ,

$$\mathbf{a}^T \mathbf{t}_n + b = 0 \quad (4.18)$$

- Then the model prediction for any value of  $\mathbf{x}$  will satisfy the same constraint

$$\mathbf{a}^T \mathbf{y}(\mathbf{x}) + b = 0 \quad (4.19)$$

- Thus, if we use 1-of- $K$  coding scheme for  $K$  classes, the elements of  $\mathbf{y}(\mathbf{x})$  will sum to 1 for any value of  $\mathbf{x}$ 
  - However, this summation constraint alone can not be interpreted as probabilities because they aren't considered to lie within the interval  $(0, 1)$



## Exercise 4.2

Consider the minimization of a sum-of-squares error function (4.15), and suppose that all of the target vectors in the training set satisfy a linear constraint

$$\mathbf{a}^T \mathbf{t}_n + b = 0 \quad (4.18)$$

where  $\mathbf{t}_n$  corresponds to the  $n^{th}$  row of the matrix  $\mathbf{T}$  in (4.15). Show that as a consequence of this constraint, the elements of the model prediction  $\mathbf{y}(\mathbf{x})$  given by the least-squares solution (4.17) also satisfy this constraint, so that

$$\mathbf{a}^T \mathbf{y}(\mathbf{x}) + b = 0 \quad (4.19)$$

To do so, assume that one of the basis functions  $\phi_0(\mathbf{x}) = 1$  so that the corresponding parameter  $\mathbf{w}_0$  plays the role of a bias.

# Decomposition of $\widetilde{\mathbf{W}}$ and $\widetilde{\mathbf{X}}$

- $\widetilde{\mathbf{W}} = \begin{pmatrix} \mathbf{w}_0^T \\ \mathbf{W} \end{pmatrix}$ :  $\mathbf{w}_0$  is the column vector of the bias weight
- $\widetilde{\mathbf{X}} = (\mathbf{1} \quad \mathbf{X})$ :  $\mathbf{1}$  is a column vector of  $N$  length

By using the definition above, the error function (4.15) is re-written,

$$E_D(\widetilde{\mathbf{W}}) = \frac{1}{2} \text{Tr} \left\{ (\widetilde{\mathbf{X}} \widetilde{\mathbf{W}} - \mathbf{T})^T (\widetilde{\mathbf{X}} \widetilde{\mathbf{W}} - \mathbf{T}) \right\} \quad (4.15)$$

$$= \frac{1}{2} \text{Tr} \left\{ (\mathbf{X} \mathbf{W} + \mathbf{1} \mathbf{w}_0^T - \mathbf{T})^T (\mathbf{X} \mathbf{W} + \mathbf{1} \mathbf{w}_0^T - \mathbf{T}) \right\} \quad (\text{C3})$$

# Solution for $\mathbf{w}_0$

Setting the derivative of (C3) with respect to  $\mathbf{w}_0$  to zero,

$$\begin{aligned}
 \frac{\partial E_D}{\partial \mathbf{w}_0} &= \frac{1}{2} \frac{\partial}{\partial \mathbf{w}_0} \text{Tr} \{ \{ (\mathbf{XW} + \mathbf{1w}_0^T - \mathbf{T})^T (\mathbf{XW} + \mathbf{1w}_0^T - \mathbf{T}) \} \} \\
 &= \frac{1}{2} \frac{\partial}{\partial \mathbf{w}_0} \{ \text{Tr} \{ (\mathbf{XW})^T \mathbf{XW} + (\mathbf{XW})^T \mathbf{1w}_0^T - (\mathbf{XW})^T \mathbf{T} \\
 &\quad + (\mathbf{1w}_0^T)^T \mathbf{XW} + (\mathbf{1w}_0^T)^T \mathbf{1w}_0^T - (\mathbf{1w}_0^T)^T \mathbf{T} - \mathbf{T}^T \mathbf{XW} - \mathbf{T}^T \mathbf{1w}_0^T + \mathbf{T}^T \mathbf{T} \} \} \\
 &= \frac{1}{2} \frac{\partial}{\partial \mathbf{w}_0} \text{Tr} \{ \mathbf{w}_0 \mathbf{1} \mathbf{1}^T \mathbf{w}_0^T + 2 \mathbf{w}_0 (\mathbf{1}^T \mathbf{XW} - \mathbf{1}^T \mathbf{T}) + (\mathbf{XW} - \mathbf{T})^T (\mathbf{XW} - \mathbf{T}) \} \\
 &= \frac{1}{2} \frac{\partial}{\partial \mathbf{w}_0} \text{Tr} \{ N \mathbf{w}_0 \mathbf{w}_0^T + 2 \mathbf{w}_0 \mathbf{1}^T (\mathbf{XW} - \mathbf{T}) + \text{const} \} \\
 &= \frac{N}{2} \frac{\partial}{\partial \mathbf{w}_0} \mathbf{w}_0^T \mathbf{w}_0 + \frac{2}{2} \frac{\partial}{\partial \mathbf{w}_0} \mathbf{w}_0 \{ \mathbf{1}^T (\mathbf{XW} - \mathbf{T}) \}^T + 0 \\
 &= N \mathbf{w}_0 + (\mathbf{XW} - \mathbf{T})^T \mathbf{1} = 0
 \end{aligned} \tag{C4}$$

$$N\mathbf{w}_0 + (\mathbf{X}\mathbf{W} - \mathbf{T})^T \mathbf{1} = 0 \quad (\text{C4})$$

$$\Downarrow$$

$$\mathbf{w}_0 = \frac{1}{N} \mathbf{T}^T \mathbf{1} - \frac{1}{N} \mathbf{W}^T \mathbf{X}^T \mathbf{1} = \bar{\mathbf{t}} - \mathbf{W}^T \bar{\mathbf{x}} \quad (\text{C5})$$

- $\bar{\mathbf{t}} = \frac{1}{N} \mathbf{T}^T \mathbf{1}$  is a mean vector of all target vectors
- $\bar{\mathbf{x}}$  is a mean vector of all feature vectors

# Solution for $\mathbf{W}$

Substituting (C5) into (C3),

$$\begin{aligned}
 E_D(\widetilde{\mathbf{W}}) &= \frac{1}{2} \text{Tr} \{ (\mathbf{X}\mathbf{W} + \mathbf{1}\mathbf{w}_0^T - \mathbf{T})^T (\mathbf{X}\mathbf{W} + \mathbf{1}\mathbf{w}_0^T - \mathbf{T}) \} \\
 &= \frac{1}{2} \text{Tr} \{ (\mathbf{X}\mathbf{W} + \mathbf{1}(\bar{\mathbf{t}} - \mathbf{W}^T \bar{\mathbf{x}})^T - \mathbf{T})^T (\mathbf{X}\mathbf{W} + \mathbf{1}(\bar{\mathbf{t}} - \mathbf{W}^T \bar{\mathbf{x}})^T - \mathbf{T}) \} \\
 &= \frac{1}{2} \text{Tr} \{ (\mathbf{X}\mathbf{W} + \mathbf{1}\bar{\mathbf{t}}^T - \mathbf{1}\bar{\mathbf{x}}^T \mathbf{W} - \mathbf{T})^T (\mathbf{X}\mathbf{W} + \mathbf{1}\bar{\mathbf{t}}^T - \mathbf{1}\bar{\mathbf{x}}^T \mathbf{W} - \mathbf{T}) \} \\
 &= \frac{1}{2} \text{Tr} \{ (\mathbf{X}\mathbf{W} + \bar{\mathbf{T}} - \bar{\mathbf{X}}\mathbf{W} - \mathbf{T})^T (\mathbf{X}\mathbf{W} + \bar{\mathbf{T}} - \bar{\mathbf{X}}\mathbf{W} - \mathbf{T}) \} \\
 &= \frac{1}{2} \text{Tr} \left\{ \{ (\mathbf{X} - \bar{\mathbf{X}})\mathbf{W} - (\mathbf{T} - \bar{\mathbf{T}}) \}^T \{ (\mathbf{X} - \bar{\mathbf{X}})\mathbf{W} - (\mathbf{T} - \bar{\mathbf{T}}) \} \right\} \quad (\text{C6})
 \end{aligned}$$

- $\bar{\mathbf{T}} = \mathbf{1}\bar{\mathbf{t}}^T$  is the matrix whose row is  $\bar{\mathbf{t}}^T$
- $\bar{\mathbf{X}} = \mathbf{1}\bar{\mathbf{x}}^T$  is the matrix whose row is  $\bar{\mathbf{x}}^T$

Setting the derivative of this (C6) with respect to  $\mathbf{W}$  to zero,

$$\begin{aligned}\frac{\partial E_D}{\partial \mathbf{W}} &= \frac{\partial}{\partial \mathbf{W}} \left\{ \frac{1}{2} \text{Tr} \left\{ \{(\mathbf{X} - \bar{\mathbf{X}})\mathbf{W} - (\mathbf{T} - \bar{\mathbf{T}})\}^T \{(\mathbf{X} - \bar{\mathbf{X}})\mathbf{W} - (\mathbf{T} - \bar{\mathbf{T}})\} \right\} \right\} \\ &= (\mathbf{X} - \bar{\mathbf{X}})^T \{(\mathbf{X} - \bar{\mathbf{X}})\mathbf{W} - (\mathbf{T} - \bar{\mathbf{T}})\} = 0\end{aligned}\quad (\text{C7})$$

$$\blacksquare \quad \frac{\partial}{\partial \mathbf{X}} \text{Tr} \{(\mathbf{A}\mathbf{X}\mathbf{B} + \mathbf{C})(\mathbf{A}\mathbf{X}\mathbf{B} + \mathbf{C})^T\} = 2\mathbf{A}^T(\mathbf{A}\mathbf{X}\mathbf{B} + \mathbf{C})\mathbf{B}^T$$

$$\begin{aligned}\mathbf{W} &= \{(\mathbf{X} - \bar{\mathbf{X}})^T(\mathbf{X} - \bar{\mathbf{X}})\}^{-1}(\mathbf{X} - \bar{\mathbf{X}})(\mathbf{T} - \bar{\mathbf{T}}) = \{\hat{\mathbf{X}}^T \hat{\mathbf{X}}\}^{-1} \hat{\mathbf{X}} \hat{\mathbf{T}} \\ &= \hat{\mathbf{X}}^\dagger \hat{\mathbf{T}}\end{aligned}\quad (\text{C8})$$

$$\blacksquare \quad \hat{\mathbf{X}} = \mathbf{X} - \bar{\mathbf{X}}$$

$$\blacksquare \quad \hat{\mathbf{T}} = \mathbf{T} - \bar{\mathbf{T}}$$

$$\blacksquare \quad \hat{\mathbf{X}}^\dagger = \{\hat{\mathbf{X}}^T \hat{\mathbf{X}}\}^{-1} \hat{\mathbf{X}}: \hat{\mathbf{X}}^\dagger \text{ is the pseudo-inverse of the matrix } \hat{\mathbf{X}}$$

## Verification using new input

Now, we suppose that  $\mathbf{x}^*$  is a new input,

$$\begin{aligned}
 \mathbf{y}(\mathbf{x}^*) &= \mathbf{W}^T \mathbf{x}^* + \mathbf{w}_0 = \mathbf{W}^T \mathbf{x}^* + \bar{\mathbf{t}} - \mathbf{W}^T \bar{\mathbf{x}} = \mathbf{W}^T (\mathbf{x}^* - \bar{\mathbf{x}}) + \bar{\mathbf{t}} \\
 &= (\hat{\mathbf{X}}^\dagger \hat{\mathbf{T}})^T (\mathbf{x}^* - \bar{\mathbf{x}}) + \bar{\mathbf{t}} \\
 &= \bar{\mathbf{t}} + \hat{\mathbf{T}}^T (\hat{\mathbf{X}}^\dagger)^T (\mathbf{x}^* - \bar{\mathbf{x}})
 \end{aligned} \tag{C9}$$

If we apply (4.18) to  $\bar{\mathbf{t}}$ , we get

$$\mathbf{a}^T \bar{\mathbf{t}} = \frac{1}{N} \mathbf{a}^T \mathbf{T}^T \mathbf{1} = -b \tag{C10}$$

Therefore, applying (C9) to (4.19), we obtain

$$\begin{aligned}
 \mathbf{a}^T \mathbf{y}(\mathbf{x}^*) &= \mathbf{a}^T \{ \bar{\mathbf{t}} + \hat{\mathbf{T}}^T (\hat{\mathbf{X}}^\dagger)^T (\mathbf{x}^* - \bar{\mathbf{x}}) \} \\
 &= \mathbf{a}^T \bar{\mathbf{t}} + \mathbf{a}^T \hat{\mathbf{T}}^T (\hat{\mathbf{X}}^\dagger)^T (\mathbf{x}^* - \bar{\mathbf{x}}) \\
 &= \mathbf{a}^T \bar{\mathbf{t}} = -b
 \end{aligned} \tag{C11}$$

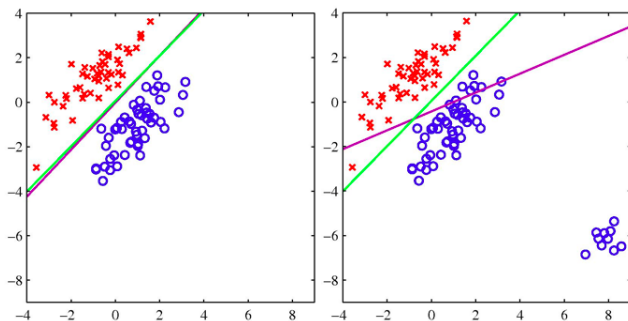
## Supplement for (C11)

$$\begin{aligned}\mathbf{a}^T \hat{\mathbf{T}}^T &= \mathbf{a}^T (\mathbf{T} - \bar{\mathbf{T}})^T = \mathbf{a}^T \mathbf{T}^T - \mathbf{a}^T \bar{\mathbf{T}}^T \\ &= \mathbf{a}^T \mathbf{T}^T - \mathbf{a}^T (\mathbf{1} \bar{\mathbf{t}}^T)^T \\ &= \mathbf{a}^T \mathbf{T}^T - \mathbf{a}^T \bar{\mathbf{t}} \mathbf{1}^T \\ &= \mathbf{a}^T \mathbf{T}^T - \frac{1}{N} \mathbf{a}^T \mathbf{T}^T \mathbf{1} \mathbf{1}^T \\ &= \mathbf{a}^T \mathbf{T}^T + b \mathbf{1}^T \\ &= -b \mathbf{1}^T + b \mathbf{1}^T \\ &= b(\mathbf{1}^T - \mathbf{1}^T) \\ &= \mathbf{0}^T\end{aligned}\tag{C12}$$



# Least Squares is highly sensitive to outliers

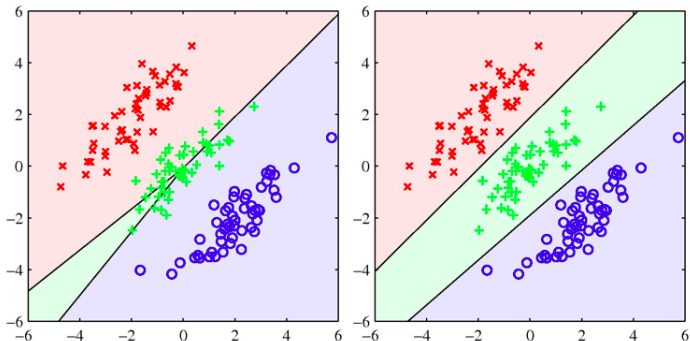
- Two-class classification between red and blue points:
  - green line: decision boundary from logistic regression
  - magenda line: decision boundary from least square



- Output range:
  - logistic regression: 0 -1
  - discriminant function in this chapter with least square: ?

# Problems with Least Squares against assumption of a Gaussian distribution

- Multiple-class classification among red, green, and blue points
  - left figure: decision boundary from least squares
  - right figure: decision boundary from logistic regression



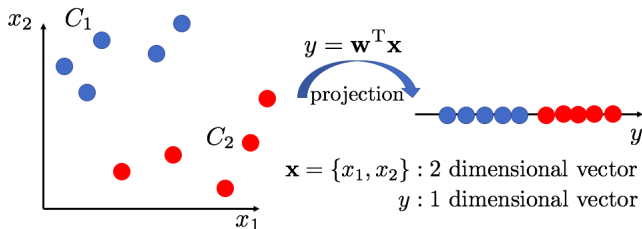
Apparently, distributions of samples aren't Gaussian distribution

# Outline

## 1 4.1 Discriminant Functions

- 4.1.1 Two classes
- 4.1.2 Multiple classes
- 4.1.3 Least squares for classification
- 4.1.4 Fisher's linear discriminant
- 4.1.5 Relation to least squares
- 4.1.6 Fisher's discriminant for multiple classes
- 4.1.7 The perceptron algorithm

# The two-class linear discriminant acts as a projection



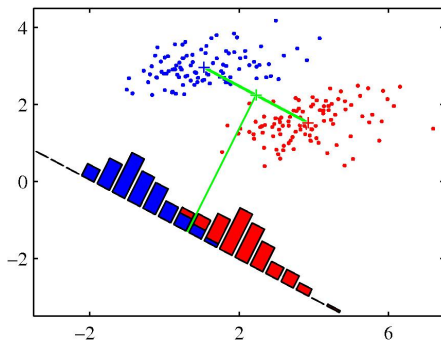
$$y = \mathbf{w}^T \mathbf{x} \quad (4.20)$$

If we place a threshold on  $y$ , this is a linear classifier

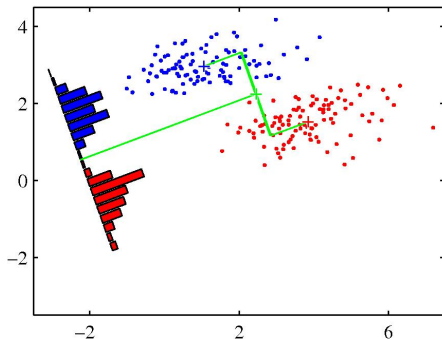
- classify  $y \geq -w_0$  as  $C_1$
- otherwise  $C_2$

# What $w$ is the best?

- A natural idea would be to project in the direction of the line connecting class means
- However, problematic if classes have variance in this direction

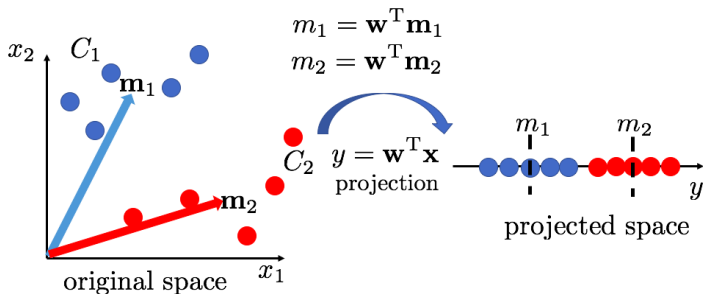


# Minimize the class overlap in the projected space



- A large separation between the projected class means
- A small variance within each class

# Mean (vector) definition in each space



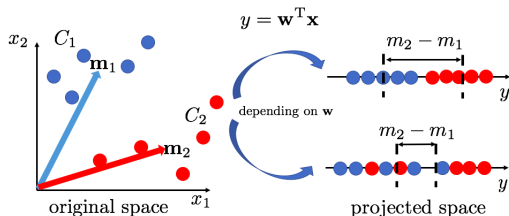
- Mean vector of each class in a original space

$$\mathbf{m}_k = \frac{1}{N_k} \sum_{n \in C_k} \mathbf{x}_n \quad (4.21^*)$$

- Mean of each class in a projected space

$$m_k = \mathbf{w}^T \mathbf{m}_k \quad (4.23)$$

# The ratio of mean separation in projected space



- The ratio of separation in projected space is distance between projected class means
  - Distance  $m_2 - m_1$  depends on the  $\mathbf{w}$

$$m_2 - m_1 = \mathbf{w}^T (\mathbf{m}_2 - \mathbf{m}_1) \quad (4.22)$$

- Now, distance  $m_2 - m_1$  can have a large value without limitation by increasing of value  $\mathbf{w}$

$\Rightarrow$  Constraint of  $\mathbf{w}$  : unit length  $\|\mathbf{w}\| = 1$



## Exercise 4.4: A Lagrange multiplier to perform the constrained maximization

- Maximize  $m_2 - m_1 = \mathbf{w}^T(\mathbf{m}_2 - \mathbf{m}_1)$
- Subject to  $\|\mathbf{w}\|^2 - 1 = 0$  ( $\|\mathbf{w}\|^2 = 1$ )
  - $\lambda$  is Lagrange multiplier

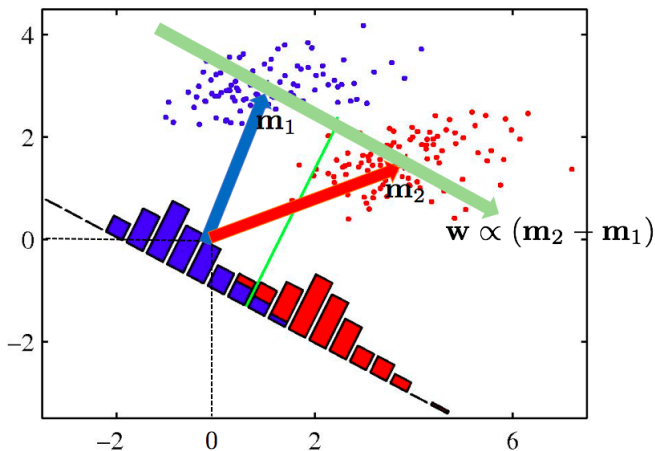
$$L(\mathbf{w}, \lambda) = \mathbf{w}^T(\mathbf{m}_2 - \mathbf{m}_1) - \lambda(\|\mathbf{w}\|^2 - 1) \quad (\text{D1})$$

- Setting the derivative of (D1) with respect to  $\mathbf{w}$  to zero

$$\frac{\partial L}{\partial \mathbf{w}} = \mathbf{m}_2 - \mathbf{m}_1 - 2\lambda\mathbf{w} = 0 \quad (\text{D2})$$

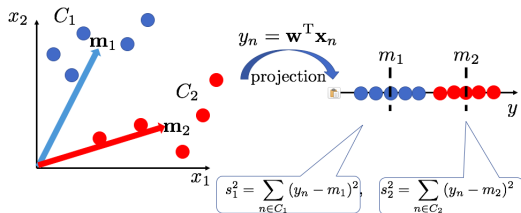
$$\mathbf{w} = \frac{1}{2\lambda}(\mathbf{m}_2 - \mathbf{m}_1) \propto (\mathbf{m}_2 - \mathbf{m}_1) \quad (\text{D3})$$

There is considerable class overlap in the projected space



$\Rightarrow$  We should take variance into consideration

# The within-class variance



- The within-class variance of the transformed (projected) data from class  $C_k$

$$s_k^2 = \sum_{n \in C_k} (y_n - m_k)^2 \quad (4.24)$$

$$y_n = \mathbf{w}^T \mathbf{x}_n \quad (4.20^*)$$

- This variance (4.24) doesn't have  $1/N$ , taking the number of samples into consideration

# Fisher criterion

$$\underline{J(\mathbf{w})} = \frac{\text{maximization : between-class variance}}{\text{minimization : within-class variance}}$$



maximization : Fisher criterion

$$J(\mathbf{w}) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2} \quad (4.25)$$

- Ratio of the between-class variance  $(m_2 - m_1)^2$  to the within-class variance  $(s_1^2 + s_2^2)$

# Deformation for maximization of $J(\mathbf{w})$

- Rewriting Fisher criterion:

$$J(\mathbf{w}) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2} = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}} \quad (4.26)$$

- Between-class covariance matrix:

$$\mathbf{S}_B = (\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^T \quad (4.27)$$

- Total within-class covariance matrix:

$$\mathbf{S}_W = \sum_{n \in C_1} (\mathbf{x}_n - \mathbf{m}_1)(\mathbf{x}_n - \mathbf{m}_1)^T + \sum_{n \in C_2} (\mathbf{x}_n - \mathbf{m}_2)(\mathbf{x}_n - \mathbf{m}_2)^T \quad (4.28)$$

# Deformation for maximization of $J(\mathbf{w})$ : Supplement

$$\mathbf{w}^T \mathbf{S}_B \mathbf{w} = \mathbf{w}^T (\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^T \mathbf{w} = (m_2 - m_1)^2 \quad (\text{D4})$$

$$\begin{aligned} \mathbf{w}^T \mathbf{S}_W \mathbf{w} &= \mathbf{w}^T \sum_{n \in C_1} (\mathbf{x}_n - \mathbf{m}_1)(\mathbf{x}_n - \mathbf{m}_1)^T \mathbf{w} \\ &\quad + \mathbf{w}^T \sum_{n \in C_2} (\mathbf{x}_n - \mathbf{m}_2)(\mathbf{x}_n - \mathbf{m}_2)^T \mathbf{w} \\ &= \sum_{n \in C_1} (y_n - m_1)^2 + \sum_{n \in C_2} (y_n - m_2)^2 \\ &= s_1^2 + s_2^2 \end{aligned} \quad (\text{D5})$$

# Maximization of $J(\mathbf{w})$

- Maximization of  $J(\mathbf{w})$ : Setting the derivative with respect  $\mathbf{w}$  to 0

$$\begin{aligned} \frac{\partial J}{\partial \mathbf{w}} &= \frac{\partial}{\partial \mathbf{w}} \left( \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}} \right) = \frac{2(\mathbf{S}_B \mathbf{w})(\mathbf{w}^T \mathbf{S}_W \mathbf{w}) - 2(\mathbf{w}^T \mathbf{S}_B \mathbf{w})(\mathbf{S}_W \mathbf{w})}{(\mathbf{w}^T \mathbf{S}_W \mathbf{w})^2} \\ &= 0 \end{aligned} \quad (\text{D6})$$

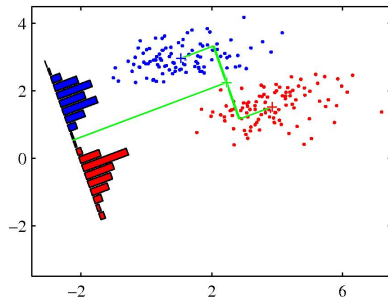
$$(\mathbf{w}^T \mathbf{S}_B \mathbf{w}) \mathbf{S}_W \mathbf{w} = (\mathbf{w}^T \mathbf{S}_W \mathbf{w}) \mathbf{S}_B \mathbf{w} \quad (4.29)$$

In term of weight vector  $\mathbf{w}$ ,

- Ignore scalar factors  $\mathbf{w}^T \mathbf{S}_B \mathbf{w}$  and  $\mathbf{w}^T \mathbf{S}_W \mathbf{w}$  because direction of weight vector  $\mathbf{w}$  only is important
- Use  $\mathbf{S}_B \mathbf{w} = (\mathbf{m}_2 - \mathbf{m}_1)((\mathbf{m}_2 - \mathbf{m}_1)^T \mathbf{w}) \propto (\mathbf{m}_2 - \mathbf{m}_1)$

$\Downarrow$

# Fisher's linear discriminant



$$\mathbf{w} \propto \mathbf{S}_W^{-1} \mathbf{S}_B \mathbf{w} \propto \mathbf{S}_W^{-1} (\mathbf{m}_2 - \mathbf{m}_1) \quad (4.30)$$

- (4.30) is a selector of direction down to a one dimension, although (4.30) is called "discriminant"
- However: we can use (4.30) to construct a classifier setting a threshold



# Outline

## 1 4.1 Discriminant Functions

- 4.1.1 Two classes
- 4.1.2 Multiple classes
- 4.1.3 Least squares for classification
- 4.1.4 Fisher's linear discriminant
- **4.1.5 Relation to least squares**
- 4.1.6 Fisher's discriminant for multiple classes
- 4.1.7 The perceptron algorithm

# Fisher criterion can be obtained as a special case of least squares especially in Two classes

- Least squares
  - Making the model predictions as close as possible to a set of target values
- Fisher criterion
  - Making maximum class separation in the projected space
- Target values in this proof (Actually not limited to this definition, Mika, "Kernel Fisher Discriminants")
  - $C_1 : \frac{N}{N_1}$
  - $C_2 : -\frac{N}{N_2}$

# Minimization of sum-of-squares error

- The sum-of-squares error function

$$E = \frac{1}{2} \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0 - t_n)^2 \quad (4.31)$$

- Setting the derivatives of  $E$  with respect to  $\mathbf{w}$  and  $w_0$  to zero

$$\frac{\partial E}{\partial \mathbf{w}} = \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0 - t_n) \mathbf{x}_n = 0 \quad (4.32)$$

$$\frac{\partial E}{\partial w_0} = \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0 - t_n) = 0 \quad (4.33)$$

- Bias  $w_0$  expression from (4.32):

$$\begin{aligned}
 \frac{\partial E}{\partial \mathbf{w}} &= \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0 - t_n) = \sum_{n=1}^N \mathbf{w}^T \mathbf{x}_n + \sum_{n=1}^N w_0 - \sum_{n=1}^N t_n \\
 &= N \mathbf{w}^T \mathbf{m} + N w_0 - \sum_{n=1}^N t_n = 0
 \end{aligned} \tag{E1}$$

$$w_0 = -\frac{N}{N} \mathbf{w}^T \mathbf{m} + \frac{1}{N} \sum_{n=1}^N t_n = -\mathbf{w}^T \mathbf{m} \tag{4.34}$$

Using following equations

$$\sum_{n=1}^N t_n = N_1 \frac{N}{N_1} - N_2 \frac{N}{N_2} = 0 \tag{4.35}$$

$$\mathbf{m} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n = \frac{1}{N} (N_1 \mathbf{m}_1 + N_2 \mathbf{m}_2) \tag{4.36}$$

- $\frac{\partial E}{\partial w_0} = 0$  (4.33) is represented by using target values that is used in this chapter

$$\left( \mathbf{S}_W + \frac{N_1 N_2}{N} \mathbf{S}_B \right) \mathbf{w} = N(\mathbf{m}_1 - \mathbf{m}_2) \quad (4.37)$$

$$(4.37) \Rightarrow \mathbf{S}_W \mathbf{w} + \frac{N_1 N_2}{N} \mathbf{S}_B \mathbf{w} = N(\mathbf{m}_1 - \mathbf{m}_2)$$

In term of  $\mathbf{w}$ ,

- Use  $\mathbf{S}_B \mathbf{w} = (\mathbf{m}_2 - \mathbf{m}_1)((\mathbf{m}_2 - \mathbf{m}_1)^T \mathbf{w}) \propto (\mathbf{m}_2 - \mathbf{m}_1)$
- Ignore scalar factors

$$\mathbf{w} \propto \mathbf{S}_W^{-1}(\mathbf{m}_2 - \mathbf{m}_1) \quad (4.38)$$

$\Downarrow$

- The discriminant function:  $y(\mathbf{x}) \geq 0 \rightarrow C_1$ , otherwise  $C_2$

$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 = \mathbf{w}^T \mathbf{x} - \mathbf{w}^T \mathbf{m} = \mathbf{w}^T (\mathbf{x} - \mathbf{m}) \quad (\text{E2})$$

# Verification from (4.33) to (4.37)

$$\begin{aligned}
 \frac{\partial E}{\partial w_0} &= \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0 - t_n) \mathbf{x}_n = \sum_{n=1}^N \mathbf{w}^T \mathbf{x}_n \mathbf{x}_n + \sum_{n=1}^N w_0 \mathbf{x}_n - \sum_{n=1}^N t_n \mathbf{x}_n \\
 &= \sum_{n=1}^N (\mathbf{x}_n \mathbf{x}_n^T) \mathbf{w} - N(\mathbf{m} \mathbf{m}^T) \mathbf{w} - N(\mathbf{m}_1 - \mathbf{m}_2) = 0
 \end{aligned} \tag{E3}$$

$$\left( \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T - N \mathbf{m} \mathbf{m}^T \right) \mathbf{w} = N(\mathbf{m}_1 - \mathbf{m}_2) \tag{E4}$$

To calculate (E4), using following equations,

$$\sum_{n=1}^N w_0 \mathbf{x}_n = N w_0 \mathbf{m} = -N(-\mathbf{w}^T \mathbf{m}) \mathbf{m} = -N(\mathbf{m} \mathbf{m}^T) \mathbf{w} \tag{E5}$$

$$\begin{aligned}
 \sum_{n=1}^N t_n \mathbf{x}_n &= \sum_{n \in C_1} t_n \mathbf{x}_n + \sum_{n \in C_2} t_n \mathbf{x}_n = N_1 \frac{N}{N_1} \mathbf{m}_1 - N_2 \frac{N}{N_2} \mathbf{m}_2 \\
 &= N(\mathbf{m}_1 - \mathbf{m}_2)
 \end{aligned} \tag{E6}$$

# Verification from (4.33) to (4.37)

- We need to prove: comparing (4.37) with (E4)

$$\sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T - N \mathbf{m} \mathbf{m}^T = \mathbf{S}_W + \frac{N_1 N_2}{N} \mathbf{S}_B \quad (\text{E7})$$

- Expanding total within-class covariance matrix

$$\begin{aligned} \mathbf{S}_W &= \sum_{n \in C_1} (\mathbf{x}_n - \mathbf{m}_1)(\mathbf{x}_n - \mathbf{m}_1)^T + \sum_{n \in C_2} (\mathbf{x}_n - \mathbf{m}_2)(\mathbf{x}_n - \mathbf{m}_2)^T \quad (4.28) \\ &= \sum_{n \in C_1} \mathbf{x}_n \mathbf{x}_n^T - 2 \sum_{n \in C_1} \mathbf{x}_n \mathbf{m}_1^T + \sum_{n \in C_1} \mathbf{m}_1 \mathbf{m}_1^T \\ &\quad + \sum_{n \in C_2} \mathbf{x}_n \mathbf{x}_n^T - 2 \sum_{n \in C_2} \mathbf{x}_n \mathbf{m}_2^T + \sum_{n \in C_2} \mathbf{m}_2 \mathbf{m}_2^T \\ &= \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T - N_1 \mathbf{m}_1 \mathbf{m}_1^T - N_2 \mathbf{m}_2 \mathbf{m}_2^T \quad (\text{E8}) \end{aligned}$$

■ Expanding the left side of the equation (E7)

$$\begin{aligned}
 \text{left} &= \mathbf{S}_W + N_1 \mathbf{m}_1 \mathbf{m}_1^T + N_2 \mathbf{m}_2 \mathbf{m}_2^T - N \mathbf{m} \mathbf{m}^T \\
 &= \mathbf{S}_W + N_1 \mathbf{m}_1 \mathbf{m}_1^T + N_2 \mathbf{m}_2 \mathbf{m}_2^T \\
 &\quad - N \left( \frac{N_1^2}{N^2} \mathbf{m}_1 \mathbf{m}_1^T + 2 \frac{N_1 N_2}{N^2} \mathbf{m}_1 \mathbf{m}_2^T + \frac{N_2^2}{N^2} \mathbf{m}_2 \mathbf{m}_2^T \right) \\
 &= \mathbf{S}_W \\
 &\quad + \frac{N_1 N_2}{N} \left( \frac{N}{N_2} \mathbf{m}_1 \mathbf{m}_1^T + \frac{N}{N_1} \mathbf{m}_2 \mathbf{m}_2^T - \frac{N_1}{N_2} \mathbf{m}_1 \mathbf{m}_1^T - 2 \mathbf{m}_1 \mathbf{m}_2^T - \frac{N_2}{N_1} \mathbf{m}_2 \mathbf{m}_2^T \right) \\
 &= \mathbf{S}_W + \frac{N_1 N_2}{N} \left( \frac{N - N_1}{N_2} \mathbf{m}_1 \mathbf{m}_1^T - 2 \mathbf{m}_1 \mathbf{m}_2^T + \frac{N - N_2}{N_1} \mathbf{m}_2 \mathbf{m}_2^T \right) \\
 &= \mathbf{S}_W + \frac{N_1 N_2}{N} (\mathbf{m}_1 \mathbf{m}_1^T - 2 \mathbf{m}_1 \mathbf{m}_2^T + \mathbf{m}_2 \mathbf{m}_2^T) \\
 &= \mathbf{S}_W + \frac{N_1 N_2}{N} (\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^T \\
 &= \mathbf{S}_W + \frac{N_1 N_2}{N} \mathbf{S}_B
 \end{aligned}$$



# Outline

## 1 4.1 Discriminant Functions

- 4.1.1 Two classes
- 4.1.2 Multiple classes
- 4.1.3 Least squares for classification
- 4.1.4 Fisher's linear discriminant
- 4.1.5 Relation to least squares
- 4.1.6 Fisher's discriminant for multiple classes
- 4.1.7 The perceptron algorithm

# Fisher's discriminant for multiple classes

- The number of classes:  $K > 2$
- Input space:  $D > K$
- Linear features:  $y_k = \mathbf{w}_k^T \mathbf{x}$ ,  $D' > 1, k = 1, \dots, D'$
- Projection matrix  $\mathbf{W}$ :  $\{\mathbf{w}_k\} = [\mathbf{w}_1 \cdots \mathbf{w}_k \cdots \mathbf{w}_{D'}]$  are the columns

$$\begin{matrix} \mathbf{y} \\ D' \times 1 \end{matrix} = \begin{matrix} \mathbf{W}^T \\ D' \times D \end{matrix} \begin{matrix} \mathbf{x} \\ D \times 1 \end{matrix}$$

- Relationship between vector  $\mathbf{y}$  and matrix  $\mathbf{W}$ :

$$\mathbf{y} = \mathbf{W}^T \mathbf{x} \quad (4.39)$$

- (4.39) transforms an original feature vector  $\mathbf{x}$  into a new feature vector  $\mathbf{y}$ , using projection matrix  $\mathbf{W}$

This (4.39) doesn't include bias.

## Definition of covariance matrix in original space

- The generalization of the total within-class covariance matrix:

$$\mathbf{S}_W = \sum_{k=1}^K \mathbf{S}_k \quad (4.40)$$

- the within-class covariance matrix:

$$\mathbf{S}_k = \sum_{n \in C_k} (\mathbf{x}_n - \mathbf{m}_k)(\mathbf{x}_n - \mathbf{m}_k)^T \quad (4.41)$$

$$\mathbf{m}_k = \frac{1}{N_k} \sum_{n \in C_k} \mathbf{x}_n \quad (4.42)$$

- Total covariance matrix:

$$\mathbf{S}_T = \sum_{n=1}^N (\mathbf{x}_n - \mathbf{m})(\mathbf{x}_n - \mathbf{m})^T \quad (4.43)$$

where  $\mathbf{m}$  is the mean of the total data set

$$\mathbf{m} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n = \frac{1}{N} \sum_{k=1}^K N_k \mathbf{m}_k \quad (4.44)$$

- Decomposition of total covariance matrix (4.43):

$$\mathbf{S}_T = \mathbf{S}_W + \mathbf{S}_B \quad (4.45)$$

where

$$\mathbf{S}_B = \sum_{k=1}^K N_k (\mathbf{m}_k - \mathbf{m})(\mathbf{m}_k - \mathbf{m})^T \quad (4.46)$$

## Definition of covariance matrix in projected space

- We can define similar matrices in the projected  $D'$ -dimensional  $\mathbf{y}$ -space

$$\mathbf{S}_W = \sum_{k=1}^K \sum_{n \in C_k} (\mathbf{y}_n - \mu_k)(\mathbf{y}_n - \mu_k)^T \quad (4.47)$$

$$\mathbf{S}_B = \sum_{k=1}^K N_k (\mu_k - \mu)(\mu_k - \mu)^T \quad (4.48)$$

where

$$\mu_k = \frac{1}{N_k} \sum_{n \in C_k} \mathbf{y}_n, \quad \mu = \frac{1}{N} \sum_{k=1}^K N_k \mu_k \quad (4.49)$$

# Construction of a scalar criterion

- Construction of a scalar criterion
  - that is large when the between-class covariance is large and when the within-class covariance is small

$$J(\mathbf{W}) = \text{Tr} \{ \mathbf{s}_W^{-1} \mathbf{s}_B \} \quad (4.50)$$

- This criterion can be rewritten as an explicit function of the projection matrix  $\mathbf{W}$

$$J(\mathbf{W}) = \text{Tr} \{ (\mathbf{W} \mathbf{S}_W \mathbf{W}^T)^{-1} (\mathbf{W} \mathbf{S}_B \mathbf{W}^T) \} \quad (4.51)$$

- The projection matrix  $\mathbf{W}$  is determined by the eigenvectors of  $\mathbf{S}_W^{-1} \mathbf{S}_B$  that correspond to the  $D'$  largest eigenvalues

$\mathbf{W}$  is determined by the eigenvectors of  $\mathbf{S}_\mathbf{W}^{-1}\mathbf{S}_\mathbf{B}$  that correspond to the  $D'$  largest eigenvalues

Taking the derivative of (4.51) with respect to  $\mathbf{W}$ ,

$$\begin{aligned}
 \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}} &= \frac{\partial}{\partial \mathbf{W}} \{ \text{Tr} \{ \mathbf{s}_\mathbf{W}^{-1} \mathbf{S}_\mathbf{B} \} \} = \frac{\partial}{\partial \mathbf{W}} \{ \text{Tr} \{ (\mathbf{W} \mathbf{S}_\mathbf{W} \mathbf{W}^T)^{-1} (\mathbf{W} \mathbf{S}_\mathbf{B} \mathbf{W}^T) \} \} \\
 &= -2 \mathbf{S}_\mathbf{W} \mathbf{W} (\mathbf{W}^T \mathbf{S}_\mathbf{W} \mathbf{W})^{-1} (\mathbf{W}^T \mathbf{S}_\mathbf{B} \mathbf{W}) (\mathbf{W}^T \mathbf{S}_\mathbf{W} \mathbf{W})^{-1} \\
 &\quad + 2 \mathbf{S}_\mathbf{B} \mathbf{W} (\mathbf{W}^T \mathbf{S}_\mathbf{W} \mathbf{W})^{-1} \\
 &= -2 \mathbf{S}_\mathbf{W} \mathbf{W} \mathbf{s}_\mathbf{W}^{-1} \mathbf{S}_\mathbf{B} \mathbf{s}_\mathbf{W}^{-1} + 2 \mathbf{S}_\mathbf{B} \mathbf{W} \mathbf{s}_\mathbf{W}^{-1}
 \end{aligned} \tag{F1}$$

Setting equation (F1) to zero,

$$\begin{aligned}
 \mathbf{S}_\mathbf{B} \mathbf{W} \mathbf{s}_\mathbf{W}^{-1} &= \mathbf{S}_\mathbf{W} \mathbf{W} \mathbf{s}_\mathbf{W}^{-1} \mathbf{S}_\mathbf{B} \mathbf{s}_\mathbf{W}^{-1} \\
 \mathbf{S}_\mathbf{B} \mathbf{W} &= \mathbf{S}_\mathbf{W} \mathbf{W} \mathbf{s}_\mathbf{W}^{-1} \mathbf{S}_\mathbf{B} \\
 (\mathbf{S}_\mathbf{W}^{-1} \mathbf{S}_\mathbf{B}) \mathbf{W} &= \mathbf{W} (\mathbf{s}_\mathbf{W}^{-1} \mathbf{S}_\mathbf{B})
 \end{aligned} \tag{F2}$$

# Preparation for proof

- Two matrices  $\mathbf{s}_B$  and  $\mathbf{s}_W$  can be simultaneously diagonalized to  $\mathbf{P}$  and  $\mathbf{Q}$  by a linear transformation  $\mathbf{z} = \mathbf{B}^T \mathbf{y}$

$$\mathbf{B}^T \mathbf{s}_B \mathbf{B} = \mathbf{P}_b, \quad \mathbf{B}^T \mathbf{s}_W \mathbf{B} = \mathbf{Q}_w \quad (\text{F3})$$

where  $\mathbf{B}$  is a  $D' \times D'$  regular matrix and  $\mathbf{B}^{-1}$  exists

- The criterion value is invariant under this regular mapping from  $\mathbf{y}$  to  $\mathbf{z}$

$$\begin{aligned} \text{Tr} \{ \mathbf{Q}_w^{-1} \mathbf{P}_b \} &= \text{Tr} \{ (\mathbf{B}^T \mathbf{s}_W \mathbf{B})^{-1} (\mathbf{B}^T \mathbf{s}_B \mathbf{B}) \} = \text{Tr} \{ \mathbf{B}^{-1} \mathbf{s}_W^{-1} \mathbf{B}^T \mathbf{B}^T \mathbf{s}_B \mathbf{B} \} \\ &= \text{Tr} \{ \mathbf{B}^{-1} \mathbf{s}_W^{-1} \mathbf{s}_B \mathbf{B} \} = \text{Tr} \{ \mathbf{s}_W^{-1} \mathbf{s}_B \} \end{aligned} \quad (\text{F4})$$

- "Similarity invariance": The traces are invariant with respect to similarity, i.e. if  $P$  is a regular matrix,  $\text{Tr}(P^{-1}XP) = \text{Tr}(X)$



Condition that satisfy (F2), to maximize the criterion (4.50)

$$(\mathbf{S}_W^{-1}\mathbf{S}_B)\mathbf{W} = \mathbf{W}(\mathbf{s}_W^{-1}\mathbf{s}_B) \quad (\text{F2})$$

(F2) is rewritten by using (F3)

$$\begin{aligned} (\mathbf{S}_W^{-1}\mathbf{S}_B)\mathbf{W} &= \mathbf{W}(\mathbf{B}\mathbf{Q}_w^{-1}\mathbf{B}^{-1}\mathbf{B}\mathbf{P}_b\mathbf{B}^{-1}) \\ (\mathbf{S}_W^{-1}\mathbf{S}_B)\mathbf{W} &= \mathbf{W}(\mathbf{B}\mathbf{Q}_w^{-1}\mathbf{P}_b\mathbf{B}^{-1}) \\ (\mathbf{S}_W^{-1}\mathbf{S}_B)\mathbf{W}\mathbf{B} &= \mathbf{W}\mathbf{B}\mathbf{Q}_w^{-1}\mathbf{P}_b \\ (\mathbf{S}_W^{-1}\mathbf{S}_B)\mathbf{W}\mathbf{B} &= (\mathbf{W}\mathbf{B})\mathbf{\Lambda} \end{aligned} \quad (\text{F5})$$

- The components of  $\mathbf{\Lambda}$  and the column vectors of  $\mathbf{W}\mathbf{B}$  are the  $D'$  eigenvalues and eigenvectors of  $\mathbf{S}_W^{-1}\mathbf{S}_B$  respectively

The trace of a matrix is the summation of the eigenvalues

- $\mu$  and  $\lambda$  are eigenvalues of  $\{\mathbf{S}_W^{-1}\mathbf{S}_B\}$  and  $\{\mathbf{s}_W^{-1}\mathbf{s}_B\}$  respectively

$$\text{Tr} \{ \mathbf{S}_W^{-1} \mathbf{S}_B \} = \mu_1 + \cdots + \mu_{D'} + \cdots + \mu_D \quad (\text{F6})$$

$$J(\mathbf{W}) = \text{Tr} \{ \mathbf{s}_W^{-1} \mathbf{s}_B \} = \lambda_1 + \cdots + \lambda_{D'} \quad (\text{F7})$$

- The eigenvalues  $\lambda$  of (F7) are also the eigenvalues of (F6), from (F5)

Therefore, the projection matrix  $\mathbf{W}$  is determined by eigenvectors of  $\mathbf{S}_W^{-1}\mathbf{S}_B$  that correspond to the  $D'$  largest eigenvalues

$K - 1$  or more future can not be found by this means

- $\mathbf{S}_B$  is composed of the sum of  $K$  matrices
  - Each of  $K$  matrices from an outer product of two vectors is rank 1
  - Only  $(K - 1)$  of these matrices are independent as a result of constraint (4.44) that is linear combination in equation of mean
- Thus,  $\mathbf{S}_B$  has rank at most equal to  $(K - 1)$  and has nonzero eigenvalues at most  $(K - 1)$



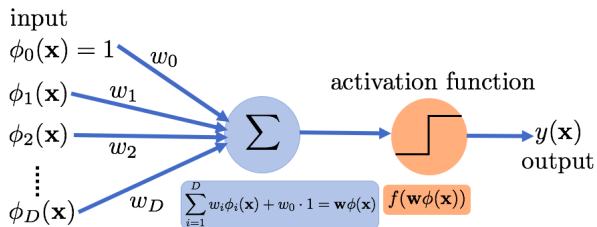
This shows that the projection onto the  $(K - 1)$ -dimensional subspace spanned by the eigenvectors of  $\mathbf{S}_B$  does not alter the value of  $J(\mathbf{W})$

# Outline

## 1 4.1 Discriminant Functions

- 4.1.1 Two classes
- 4.1.2 Multiple classes
- 4.1.3 Least squares for classification
- 4.1.4 Fisher's linear discriminant
- 4.1.5 Relation to least squares
- 4.1.6 Fisher's discriminant for multiple classes
- 4.1.7 The perceptron algorithm

# Perceptron: Two class model



- A generalized linear model:

$$y(\mathbf{x}) = f(\mathbf{w}^T \phi(\mathbf{x})) \quad (4.52)$$

- The nonlinear activation function:

$$f(a) = \begin{cases} +1, & a \geq 0 \\ -1, & a < 0 \end{cases} \quad (4.53)$$

- The main difference compared to the methods we've seen so far is the learning algorithm

# Perceptron criterion

- Target values:  $t = 1$  for class  $C_1$  and  $t = -1$  for class  $C_2$
- Minimization of the number of mis-classified examples
  - An example is mis-classified if  $\mathbf{w}^T \phi(\mathbf{x}_n) t_n < 0$
  - Perceptron criterion:

$$E_p(\mathbf{w}) = - \sum_{n \in \mathcal{M}} \mathbf{w}^T \phi(\mathbf{x}_n) t_n \quad (4.54)$$

Sum over mis-classified examples only

# Perceptron Learning Algorithm

- Minimize the error function using stochastic gradient descent (gradient descent per example):

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_P(\mathbf{w}) = \mathbf{w}^{(\tau)} + \eta \phi(\mathbf{x}_n) t_n \quad (4.55)$$

- Iterate over all training examples, only change  $\mathbf{w}$  if the example is mis-classified
- Guarantee to converge only if data are linearly separable
- May take many iterations
- Sensitive to initialization

# The effect of a single update in the perceptron learning algorithm

$$-\mathbf{w}^{(\tau+1)\text{T}}\phi(\mathbf{x}_n)t_n = -\mathbf{w}^{(\tau)\text{T}}\phi(\mathbf{x}_n)t_n - (\phi(\mathbf{x}_n)t_n)^{\text{T}}\phi(\mathbf{x}_n)t_n < -\mathbf{w}^{(\tau)\text{T}}\phi(\mathbf{x}_n)t_n \quad (4.56)$$

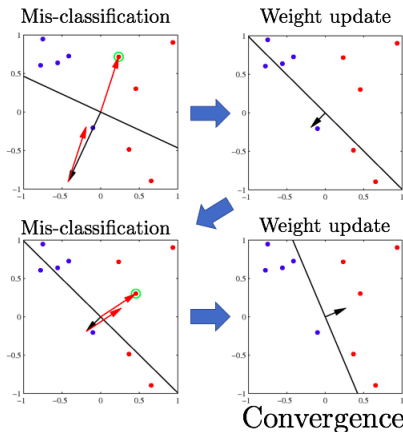
- The contribution to the error from a mis-classified pattern decrease
- The contribution to the error from the other mis-classified patterns may not decrease



The perceptron learning rule is not guaranteed to reduce the total error function at each stage



# The convergence of the perceptron learning algorithm



- Note that there are many hyperplanes with 0 error
  - Support vector machine is a nice way of choosing one