

Machine Learning

Lecture 06

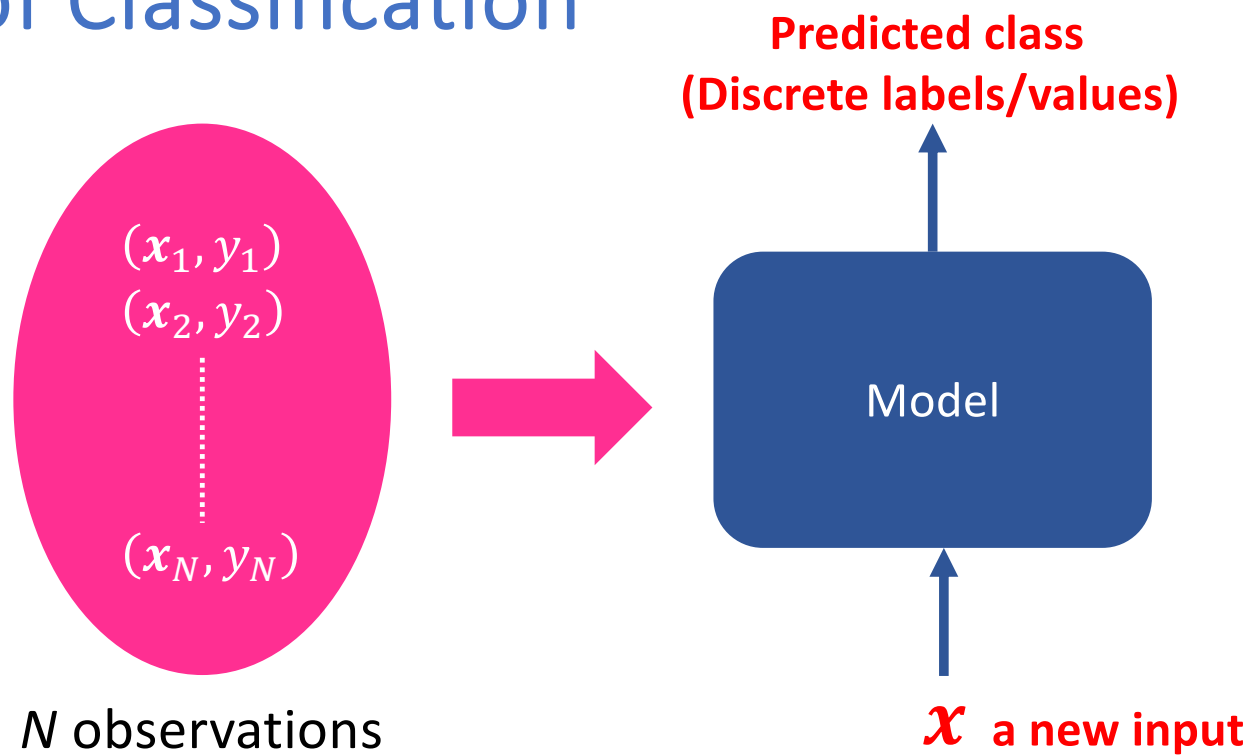
Min-Kuan Chang
minkuanc@nchu.edu.tw
EE, College of EECS

Linear Models for Classification

Topic 01

Two-Class Classification

Goal of Classification



- The goal in classification is to take an input vector \mathbf{x} and to assign it to one of K discrete classes C_k where $k = 1, 2, \dots, K$
 - the classes are taken to be disjoint
 - each input is assigned to one and only one class
 - the input space is divided into decision regions whose boundaries are called decision boundaries or decision surfaces
- In this lecture, we consider linear models for classification
 - the decision surfaces are linear functions of the input vector
 - the decision surfaces are defined by $(D - 1)$ -dimensional hyperplanes within the D -dimensional input space
 - data sets whose classes can be separated exactly by linear decision surfaces are said to be linearly separable

Two-Class Classification

- In the two-class classification, a target has two possible labels or values. For example, $y_i \in \{C_1, C_2\}$ or $y_i \in \{-1, 1\}$
- A discriminant is present to help classify an input
- The discriminant is to map an input to one of the classes, which is either C_1 or C_2 (-1 or 1) in the case of the two-class classification
- The simplest form of discriminant is the linear discriminant function

判別

$$\hat{y}(x) = w^T x + w_0$$

together with the decision boundary to assign a class to the input x

Two-Class Classification

- In the two-class classification, the linear discriminant function maps an input \mathbf{x} to one of the two classes, C_1 and C_2

$$\hat{y}(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 \begin{matrix} \overset{C_1}{\geq} \\ \underset{C_2}{\leq} \end{matrix} 0$$

- The decision boundary in the two-class classification

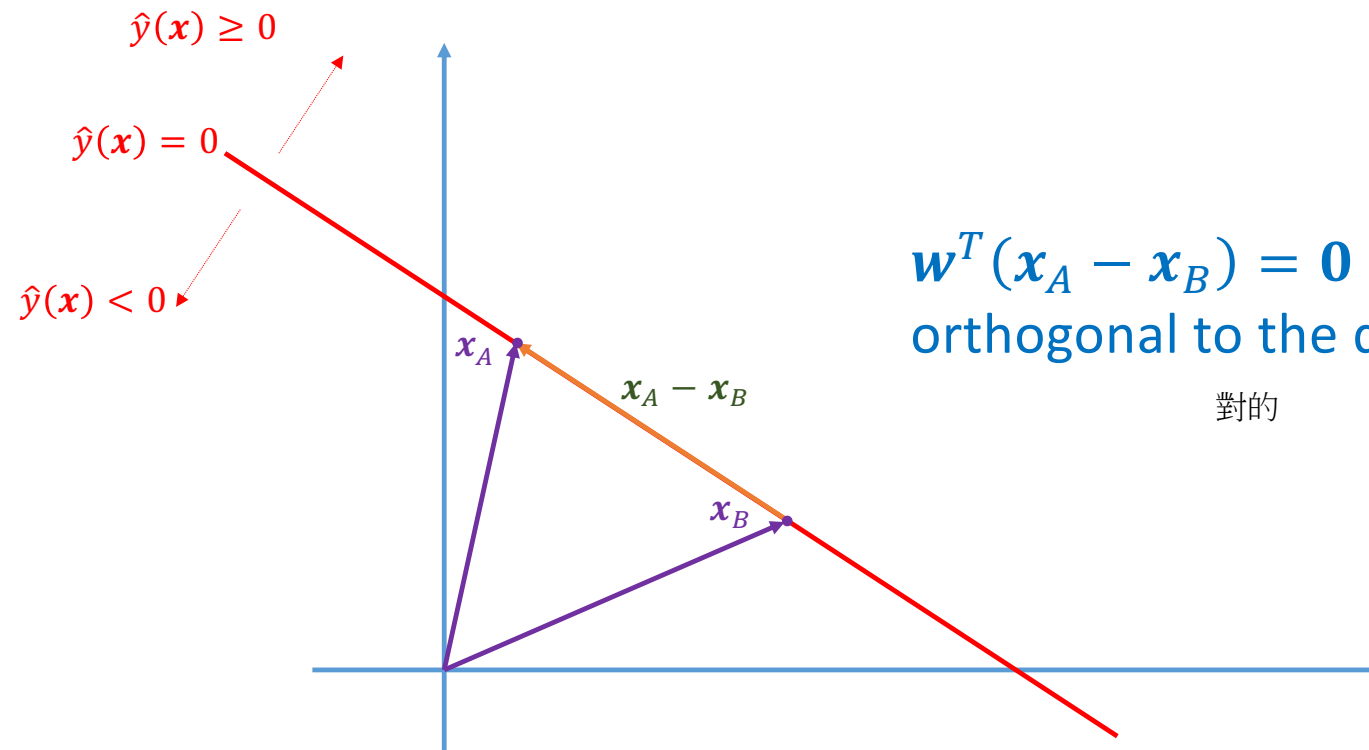
$$\hat{y}(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 = 0$$

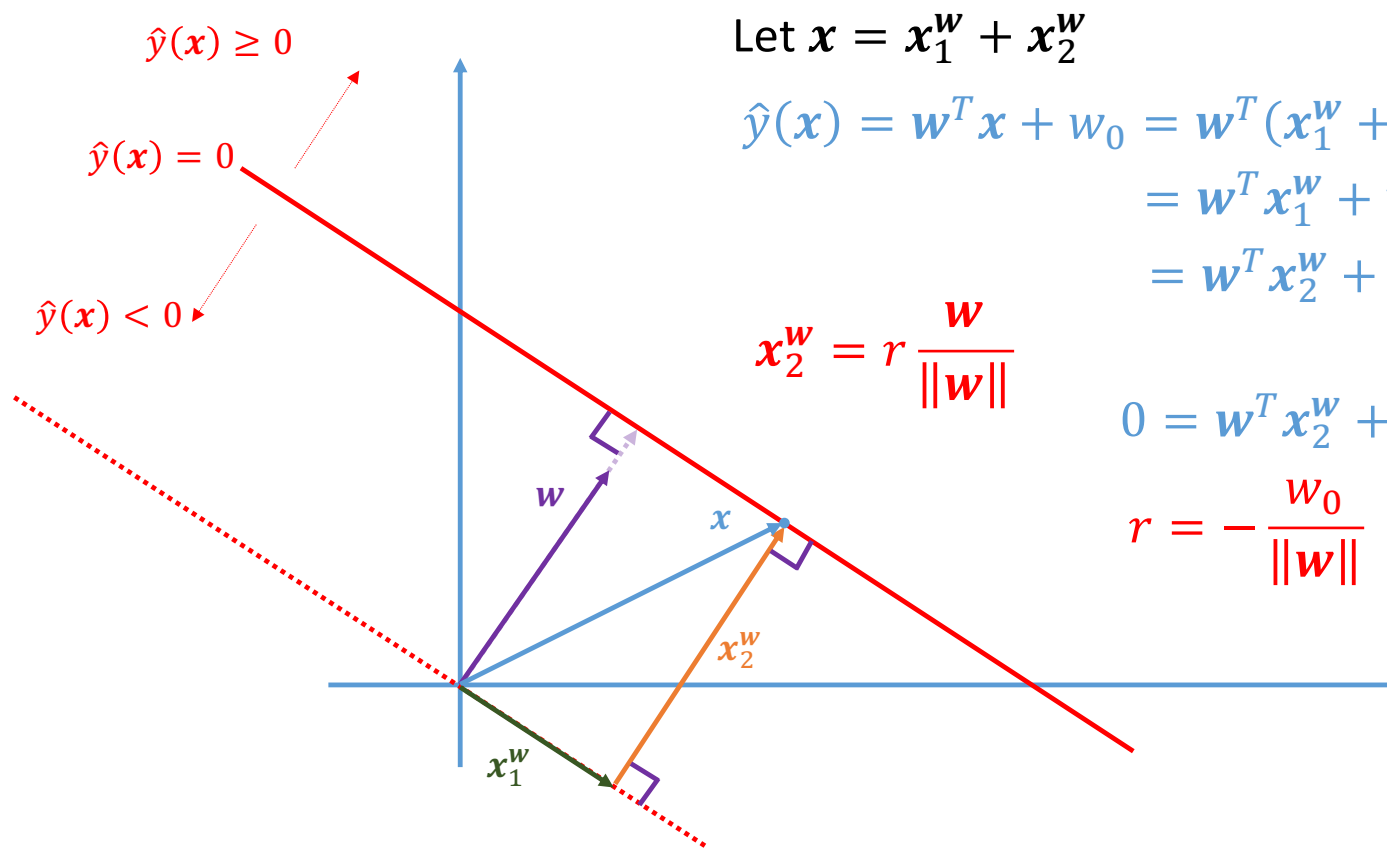
Two-Class Classification

- Geometry of the decision boundary
 - two inputs, say x_A and x_B lie on the decision boundary

The diagram illustrates the derivation of the decision boundary equation. At the top, two equations are shown: $w^T x_A + w_0 = 0$ on the left and $w^T x_B + w_0 = 0$ on the right. Arrows from each equation point towards a central circle containing a minus sign. From the bottom of this circle, an arrow points down to the final equation: $w^T (x_A - x_B) = 0$.

$$w^T x_A + w_0 = 0 \qquad w^T x_B + w_0 = 0$$
$$w^T (x_A - x_B) = 0$$





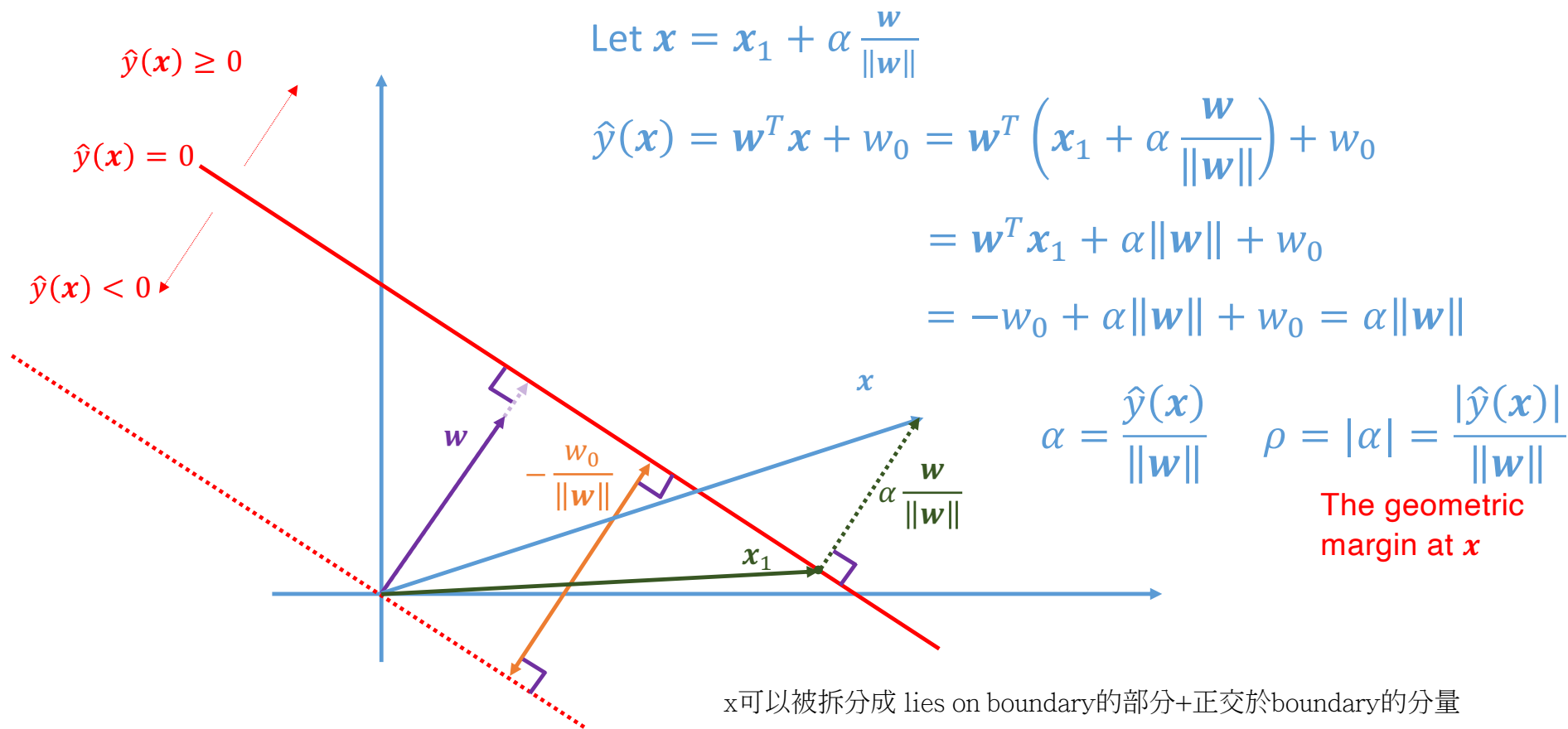
$$\text{Let } x = x_1^w + x_2^w$$

$$\begin{aligned}\hat{y}(x) &= w^T x + w_0 = w^T (x_1^w + x_2^w) + w_0 \\ &= w^T x_1^w + w^T x_2^w + w_0 \\ &= w^T x_2^w + w_0 = 0\end{aligned}$$

$$x_2^w = r \frac{w}{\|w\|}$$

$$0 = w^T x_2^w + w_0 = r \|w\| + w_0$$

$$r = -\frac{w_0}{\|w\|}$$



Multi-class Classification

- The goal is to assign an input to one of the K classes
- Adopting the binary classification leads to the ambiguity
 - an one-versus-the-rest classifier
 - an one-versus-one classifier



Multi-class Classification

- To avoid the ambiguity, a single K -class discriminant comprising K linear functions of the form

$$\hat{y}_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k,0}$$

for $k = 1, 2, \dots, K$

- An input \mathbf{x} is assigned to class j if

$$\hat{y}_j(\mathbf{x}) > \hat{y}_k(\mathbf{x}) \quad \text{or} \quad j = \arg \max_{k=1,2,\dots,K} \hat{y}_k(\mathbf{x})$$

for $k = 1, 2, \dots, K$ and $k \neq j$

Linear Models for Classification

Topic 02

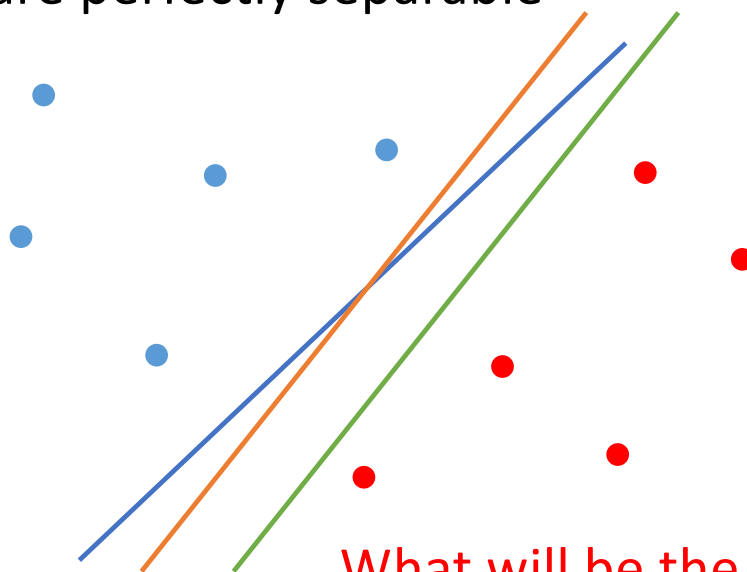
Support Vector Machine Separable Case - Part I

Linear model

Nonlinear model => "kernel"

The Separable Case

- We first assume that separable case in which the collected points of the two classes are perfectly separable



What will be the best decision boundary?

The Separable Case

- Recall the geometric margin,

$$\rho(\mathbf{x}) = \frac{|\hat{y}(\mathbf{x})|}{\|\mathbf{w}\|}$$

- The minimum geometric margin of $\hat{y}(\mathbf{x})$ given \mathbf{w}

$$\rho = \min_{n \in \{1, 2, \dots, N\}} \frac{|\hat{y}(\mathbf{x}_n)|}{\|\mathbf{w}\|}$$

where \mathbf{x}_n for $n = 1, 2, \dots, N$ are the input points

[註解]

1. 模型：

$$= \hat{y}(\mathbf{x}) = \mathbf{w}\mathbf{x} + w_0,$$

是一個分隔面，將資料切分成兩區的平面。

2. minimum geometric margin：

在 n 個 samples 中，找到一個最小化 $\hat{y}(\mathbf{x})$ 的 sample，
白話一點就是找到一個 sample 離 \hat{y} 平面最近，
其計算出的值除以 w 的長度，
就是 minimum geometric margin。

The Separable Case

- The safest way is to adjust \mathbf{w} and w_0 so that ρ can be as large as possible
我們的主要目標是，將分隔平面(\hat{y})離最近的資料越遠越好
- Given the \mathbf{x}_n for $n = 1, 2, \dots, N$ and its corresponding target $y_i \in \{-1, 1\}$, the maximum-margin optimization is

$$\rho_{max} = \max_{\mathbf{w}, w_0} \rho = \max_{\mathbf{w}, w_0} \min_n \frac{|\hat{y}(\mathbf{x}_n)|}{\|\mathbf{w}\|} = \max_{\mathbf{w}, w_0} \min_n \frac{|\mathbf{w}^T \mathbf{x}_n + w_0|}{\|\mathbf{w}\|}$$

subject to

因此我們藉由調整 \mathbf{w} , w_0 ，來最大化「離分隔面最近的資料」到「分隔面」的距離

$$y_n \hat{y}(\mathbf{x}_n) = y_n (\mathbf{w}^T \mathbf{x}_n + w_0) \geq 0$$

屬於 class = +1 類的資料點 \mathbf{x} ，其 $\hat{y}(\mathbf{x}_n) > 0$ ，其類別 $y_n = +1$
屬於 class = -1 類的資料點則 $\hat{y}(\mathbf{x}) < 0$ ，類別 $y_n = -1$ ，

因此相乘必為正值，落在分隔面上則 = 0。

The Separable Case

- Note that $y_n(\mathbf{w}^T \mathbf{x}_n + w_0) \geq 0$ and $y_n \in \{-1, 1\}$. This implies that

$$|\mathbf{w}^T \mathbf{x}_n + w_0| = y_n(\mathbf{w}^T \mathbf{x}_n + w_0)$$

所以我們可以拿掉絕對值，用「所屬類別」乘以「模型預測值」來代替。

- The maximum-margin optimization becomes

$$\rho_{max} = \max_{\mathbf{w}, w_0} \min_n \frac{y_n(\mathbf{w}^T \mathbf{x}_n + w_0)}{\|\mathbf{w}\|}$$

The Separable Case

- Observing

- when we make the rescaling $\mathbf{w} \rightarrow \kappa \mathbf{w}$ and $w_0 \rightarrow \kappa w_0$ for $\kappa > 0$

$$\frac{y_n(\mathbf{w}^T \mathbf{x}_n + w_0)}{\|\mathbf{w}\|} = \frac{\kappa y_n(\mathbf{w}^T \mathbf{x}_n + w_0)}{\kappa \|\mathbf{w}\|} = \frac{y_n(\mathbf{w}^T \mathbf{x}_n + w_0)}{\|\mathbf{w}\|}$$

- we can rescale \mathbf{w} and w_0 such that

$$\min_{n=1,2,\dots,N} y_n(\mathbf{w}^T \mathbf{x}_n + w_0) = 1$$

我們習慣將離分隔面最近的點之預測值設(rescale)為1

The Separable Case

- Based on this observation, the maximum-margin optimization can be rewritten as

$$\max_{\mathbf{w}, w_0} \frac{1}{\|\mathbf{w}\|}$$

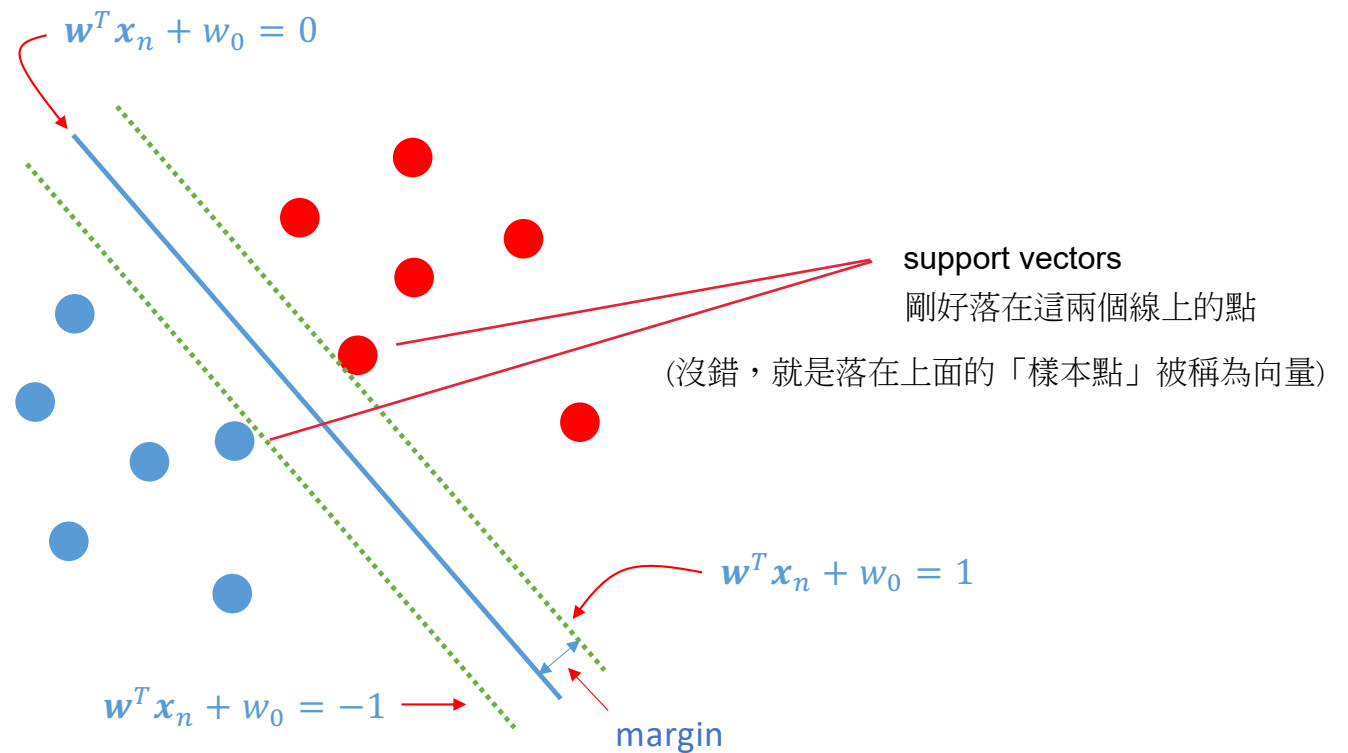
由於最近的點的預測值的絕對值被定為1，
分子部分就能寫成1，我們僅需處理分母 \mathbf{w} 的部分。

such that $y_n(\mathbf{w}^T \mathbf{x}_n + w_0) \geq 1$ for $n = 1, 2, \dots, N$. This optimization problem is equivalent to

$$\min_{\mathbf{w}, w_0} \frac{1}{2} \|\mathbf{w}\|^2$$

subject to $y_n(\mathbf{w}^T \mathbf{x}_n + w_0) \geq 1$

The Separable Case



The Separable Case

- This constraint optimization

$$\min_{\mathbf{w}, w_0} \frac{1}{2} \|\mathbf{w}\|^2 \quad \text{subject to } y_n(\mathbf{w}^T \mathbf{x}_n + w_0) \geq 1$$

can be converted to

$$\min_{\mathbf{w}, w_0} J(\mathbf{w}, w_0, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{n=1}^N \alpha_n (1 - y_n(\mathbf{w}^T \mathbf{x}_n + w_0))$$

where $\boldsymbol{\alpha} = [\alpha_1, \alpha_2, \dots, \alpha_N]^T$ and $\alpha_n \geq 0$ for $n = 1, 2, \dots, N$

The Separable Case

- Applying KKT conditions at the optimum, we have

- $\frac{\partial}{\partial \mathbf{w}} J(\mathbf{w}, w_0, \alpha) = 0$
求 \mathbf{w}

$$\mathbf{w} - \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n = 0$$

此式很重要，後續會用來代換 \mathbf{w}

- $\frac{\partial}{\partial w_0} J(\mathbf{w}, w_0, \alpha) = 0$
求 w_0

$$\sum_{n=1}^N -\alpha_n y_n = 0$$

- complementary slackness

$$\alpha_n (1 - y_n (\mathbf{w}^T \mathbf{x}_n + w_0)) = 0$$

The Separable Case

- Discussion:
 - the complementary slackness condition:
 - when $\alpha_n = 0$, $1 - y_n(\mathbf{w}^T \mathbf{x}_n + w_0) > 0$
 - when $\alpha_n > 0$, $1 - y_n(\mathbf{w}^T \mathbf{x}_n + w_0) = 0$
 - the complementary slackness condition together with $\mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n$:
 - when $\alpha_n > 0$, \mathbf{x}_n will contribute to \mathbf{w} and we call such \mathbf{x}_n the support vector
 - when \mathbf{x}_n is the support vector, $1 - y_n(\mathbf{w}^T \mathbf{x}_n + w_0) = 0$ and \mathbf{x}_n lies on the marginal hyperplanes, $\mathbf{w}^T \mathbf{x}_n + w_0 = 1$ or $\mathbf{w}^T \mathbf{x}_n + w_0 = -1$
 - $\min_{\mathbf{w}, w_0} J(\mathbf{w}, w_0, \boldsymbol{\alpha})$ is a strict convex optimization problem of \mathbf{w}
 - $\frac{\partial}{\partial \mathbf{w} \partial \mathbf{w}} J(\mathbf{w}, w_0, \boldsymbol{\alpha})$ is positive definite
 - $J(\mathbf{w}, w_0, \boldsymbol{\alpha})$ is strictly convex
 - the optimal \mathbf{w} is unique

Linear Models for Classification

Topic 03

Support Vector Machine Separable Case - Part II

The Separable Case

Primal Problem

- Recall

$$\min_{\mathbf{w}, w_0} J(\mathbf{w}, w_0, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{n=1}^N \alpha_n (1 - y_n (\mathbf{w}^T \mathbf{x}_n + w_0))$$

- KKT conditions

$$\mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n \quad \sum_{n=1}^N \alpha_n y_n = 0 \quad \alpha_n (1 - y_n (\mathbf{w}^T \mathbf{x}_n + w_0)) = 0 \quad \forall n \in \{1, 2, \dots, N\}$$

The Separable Case

- The dual problem

$$\max_{\alpha} \left\{ \min_{\mathbf{w}, w_0} J(\mathbf{w}, w_0, \alpha) \right\}$$

subject to $\sum_{n=1}^N \alpha_n y_n = 0$ and $\alpha_n \geq 0$ for $n = 1, 2, \dots, N$

The Separable Case

$$\begin{aligned}
 \min_{\mathbf{w}, w_0} J(\mathbf{w}, w_0, \alpha) &= \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{n=1}^N \alpha_n (1 - y_n (\mathbf{w}^T \mathbf{x}_n + w_0)) \\
 &= \frac{1}{2} \left\| \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n \right\|^2 + \sum_{n=1}^N \alpha_n - \sum_{n=1}^N \alpha_n y_n (\mathbf{w}^T \mathbf{x}_n + w_0) \\
 &= \frac{1}{2} \left\| \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n \right\|^2 + \sum_{n=1}^N \alpha_n - \sum_{n=1}^N \alpha_n y_n \left(\sum_{n=1}^N \alpha_n y_n \mathbf{x}_n \right)^T \mathbf{x}_n - w_0 \sum_{n=1}^N \alpha_n y_n \\
 &= \frac{1}{2} \sum_{n=1}^N \sum_{n'=1}^N \alpha_n y_n \alpha_{n'} y_{n'} \mathbf{x}_{n'}^T \mathbf{x}_n + \sum_{n=1}^N \alpha_n - \sum_{n=1}^N \sum_{n'=1}^N \alpha_n y_n \alpha_{n'} y_{n'} \mathbf{x}_{n'}^T \mathbf{x}_n \\
 &= -\frac{1}{2} \sum_{n=1}^N \sum_{n'=1}^N \alpha_n y_n \alpha_{n'} y_{n'} \mathbf{x}_{n'}^T \mathbf{x}_n + \sum_{n=1}^N \alpha_n
 \end{aligned}$$

The Separable Case

Dual Problem

- The dual problem becomes

$$\max_{\alpha} -\frac{1}{2} \sum_{n=1}^N \sum_{n'=1}^N \alpha_n y_n \alpha_{n'} y_{n'} \mathbf{x}_{n'}^T \mathbf{x}_n + \sum_{n=1}^N \alpha_n$$

subject to $\sum_{n=1}^N \alpha_n y_n = 0$ and $\alpha_n \geq 0$ for $n = 1, 2, \dots, N$

This is a quadratic programming

The Separable Case

- Let $\alpha^* = [\alpha_1^*, \alpha_2^*, \dots, \alpha_N^*]$ be the solution to the dual problem. The class assigned to an input \mathbf{x} is

$$\text{sgn}(\hat{y}(\mathbf{x})) = \text{sgn} \left(\left(\sum_{n=1}^N \alpha_n^* y_n \mathbf{x}_n \right)^T \mathbf{x} + w_0 \right)$$

- w_0 can be obtained from any support vector via

$$w_0 = \frac{1}{y_i} - \left(\sum_{n=1}^N \alpha_n^* y_n \mathbf{x}_n \right)^T \mathbf{x}_i = y_i - \left(\sum_{n=1}^N \alpha_n^* y_n \mathbf{x}_n \right)^T \mathbf{x}_i$$

The Separable Case

- The fact that

$$w_0 = \frac{1}{y_i} - \left(\sum_{n=1}^N \alpha_n^* y_n \mathbf{x}_n \right)^T \mathbf{x}_i = y_i - \left(\sum_{n=1}^N \alpha_n^* y_n \mathbf{x}_n \right)^T \mathbf{x}_i$$

gives an interesting result

$$\rho^2 = \frac{1}{\|\alpha^*\|_1}$$

which can be obtained by multiplying both sides by $\alpha_i^* y_i$ for $\alpha_i^* > 0$ and taking the sum

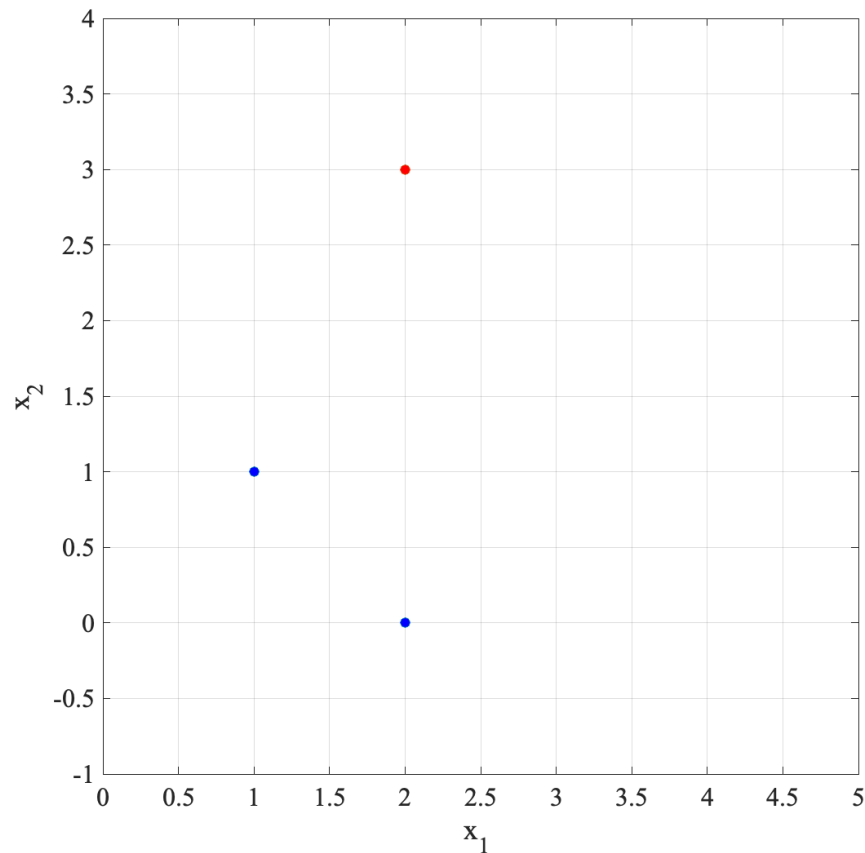
Example

- Class 1:

$$\mathbf{x}_1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \mathbf{x}_2 = \begin{bmatrix} 2 \\ 0 \end{bmatrix}$$

- Class 2:

$$\mathbf{x}_3 = \begin{bmatrix} 2 \\ 3 \end{bmatrix}$$

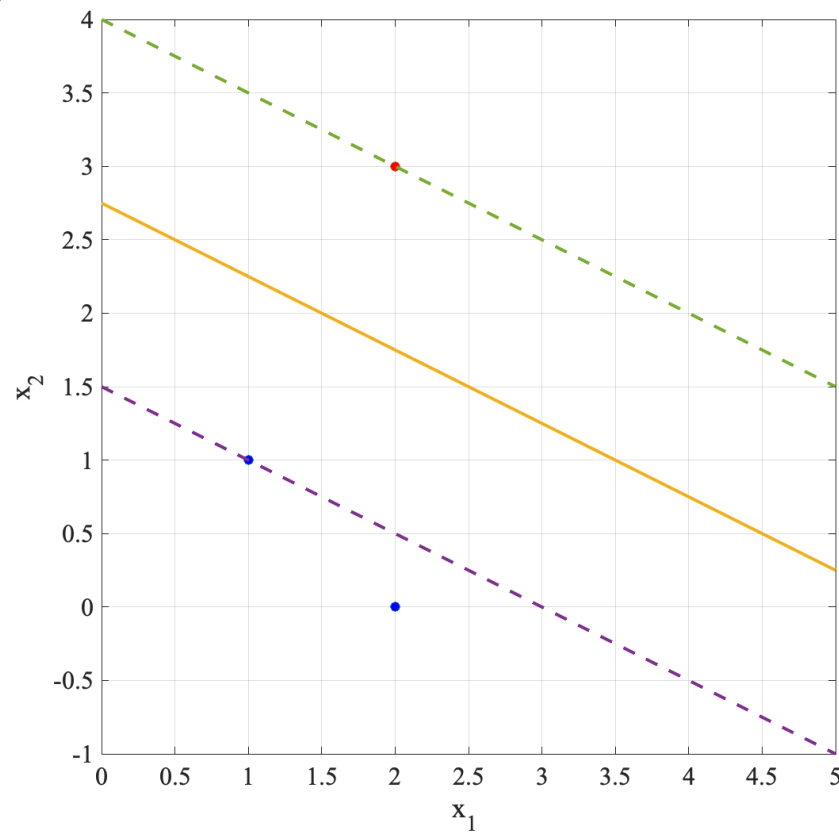


Example

- From the dual problem, using the algorithm for quadratic programming, we have $\alpha_1 = 0.8, \alpha_2 = 0, \alpha_3 = 0.8$
- $\begin{bmatrix} 1 \\ 1 \end{bmatrix}$ and $\begin{bmatrix} 2 \\ 3 \end{bmatrix}$ are support vectors
- This leads to $w_1 = 0.8$ and $w_2 = 1.6$. This indicates that $w_1 : w_2 = 1 : 2$. We let $\mathbf{w} = [a, 2a]^T$
- We know that

$$w^T \begin{bmatrix} 1 \\ 1 \end{bmatrix} + w_0 = -1 \text{ and } w^T \begin{bmatrix} 2 \\ 3 \end{bmatrix} + w_0 = 1 \quad \Rightarrow \quad a = 0.5, w_0 = -2.2$$

Example



$$\hat{y} \left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \right) = 0.4x_1 + 0.8x_2 - 2.2$$

Linear Models for Classification

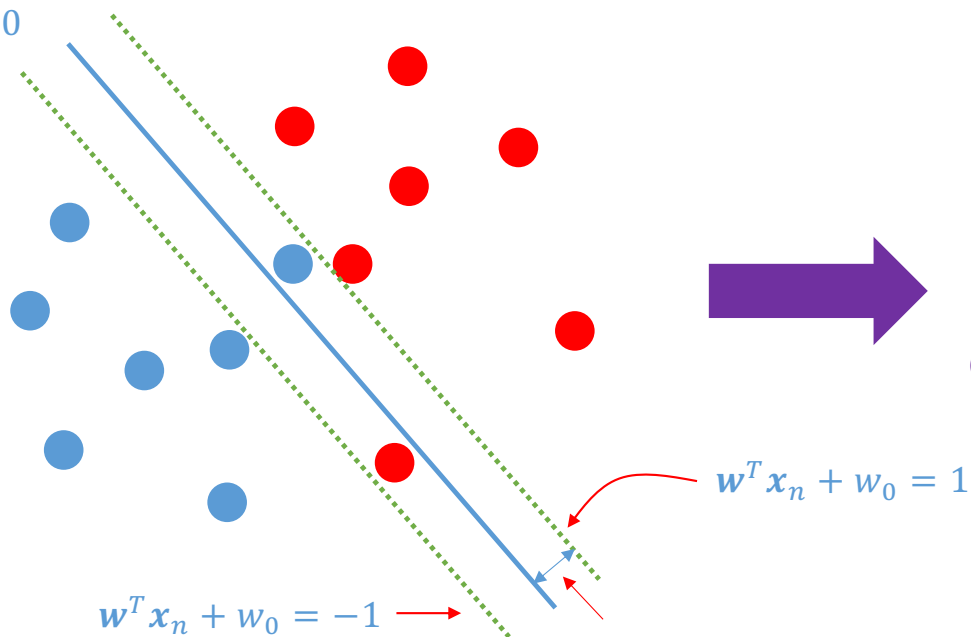
Topic 04

Support Vector Machine Non-Separable Case

The Non-Separable Case

- In many cases, points of different classes are not always separable

$$\mathbf{w}^T \mathbf{x}_n + w_0 = 0$$



$$y_n(\mathbf{w}^T \mathbf{x}_n + w_0) \not\geq 1$$

Condition for being separable is violated

The Non-Separable Case

- Introducing the slack variables to relax the optimization problem in the separable case

$$y_n(\mathbf{w}^T \mathbf{x}_n + w_0) \geq 1 - \epsilon_n$$

where $\epsilon_n \geq 0$

- ϵ_n measures the distance by which \mathbf{x}_n violates $y_n(\mathbf{w}^T \mathbf{x}_n + w_0) \geq 1$
- \mathbf{x}_n with $\epsilon_n > 0$ is considered as the outlier
- those outlier with $0 < \epsilon_n < 1$ can be still be classified correctly.
- for this reason, the margin ρ is called the soft margin as opposed to the hard margin in the separable case

The Non-Separable Case

- In the non-separable case, two conflicting objects are faced
 - we would like to limit the L_p -norm of these slack variables
 - results in fewer outliers
 - leads to the small margin
 - we would like to find the decision boundary with the margin as large as possible
 - causes more outliers
 - leads to the large L_p -norm of these slack variables

The Non-Separable Case

- Based on these two conflicting goals, the optimization becomes

$$\min_{\mathbf{w}, w_0} \frac{1}{2} \|\mathbf{w}\|^2 + \gamma \sum_{n=1}^N \epsilon_n^q$$

subject to

$$y_n(\mathbf{w}^T \mathbf{x}_n + w_0) \geq 1 - \epsilon_n, n = 1, 2, \dots, N$$

$$\epsilon_n \geq 0, n = 1, 2, \dots, N$$

The Non-Separable Case

- Introducing the Lagrange multipliers, we can convert the constrained optimization into unconstrained one
- When $q = 1$, The objective function becomes

$$J(\mathbf{w}, w_0, \boldsymbol{\epsilon}, \boldsymbol{\alpha}, \boldsymbol{\beta}) = \frac{1}{2} \|\mathbf{w}\|^2 + \gamma \sum_{n=1}^N \epsilon_n + \sum_{n=1}^N \alpha_n (1 - \epsilon_n - y_n(\mathbf{w}^T \mathbf{x}_n + w_0)) - \sum_{n=1}^N \beta_n \epsilon_n$$

where $\boldsymbol{\epsilon} = [\epsilon_1, \dots, \epsilon_N]$, $\boldsymbol{\alpha} = [\alpha_1, \dots, \alpha_N]$, $\boldsymbol{\beta} = [\beta_1, \dots, \beta_N]$

$$\alpha_i \geq 0, \beta_i \geq 0, i = 1, 2, \dots, N$$

Primal Problem

The Non-Separable Case

- Applying KKT conditions at the optimum, we have

- $\frac{\partial}{\partial \mathbf{w}} J(\mathbf{w}, w_0, \epsilon, \alpha, \beta)$

$$\mathbf{w} - \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n = 0$$

- $\frac{\partial}{\partial w_0} J(\mathbf{w}, w_0, \epsilon, \alpha, \beta)$

$$\sum_{n=1}^N -\alpha_n y_n = 0$$

- $\frac{\partial}{\partial \epsilon_i} J(\mathbf{w}, w_0, \epsilon, \alpha, \beta)$

$$\gamma - \alpha_i - \beta_i = 0$$

The Non-Separable Case

- Applying KKT conditions at the optimum, we have
 - complementary slackness

$$\alpha_n(1 - \epsilon_n - y_n(\mathbf{w}^T \mathbf{x}_n + w_0)) = 0, n = 1, 2, \dots, N$$

$$\beta_n \epsilon_n = 0, n = 1, 2, \dots, N$$

The Non-Separable Case

Dual Problem

- The dual problem becomes

$$\max_{\alpha} -\frac{1}{2} \sum_{n=1}^N \sum_{n'=1}^N \alpha_n y_n \alpha_{n'} y_{n'} \mathbf{x}_{n'}^T \mathbf{x}_n + \sum_{n=1}^N \alpha_n$$

subject to $\sum_{n=1}^N \alpha_n y_n = 0$ and $0 \leq \alpha_n \leq \gamma$ for $n = 1, 2, \dots, N$

This is the same as the separable case!!
(They only differ in the constraint)

The Non-Separable Case

- Let $\alpha^* = [\alpha_1^*, \alpha_2^*, \dots, \alpha_N^*]$ be the solution to the dual problem. The class assigned to an input \mathbf{x} is

$$\text{sgn}(\hat{y}(\mathbf{x})) = \text{sgn} \left(\left(\sum_{n=1}^N \alpha_n^* y_n \mathbf{x}_n \right)^T \mathbf{x} + w_0 \right)$$

- w_0 can be obtained from any support vector **lying on a marginal hyperplane** via

$$w_0 = \frac{1}{y_i} - \left(\sum_{n=1}^N \alpha_n^* y_n \mathbf{x}_n \right)^T \mathbf{x}_i = y_i - \left(\sum_{n=1}^N \alpha_n^* y_n \mathbf{x}_n \right)^T \mathbf{x}_i$$

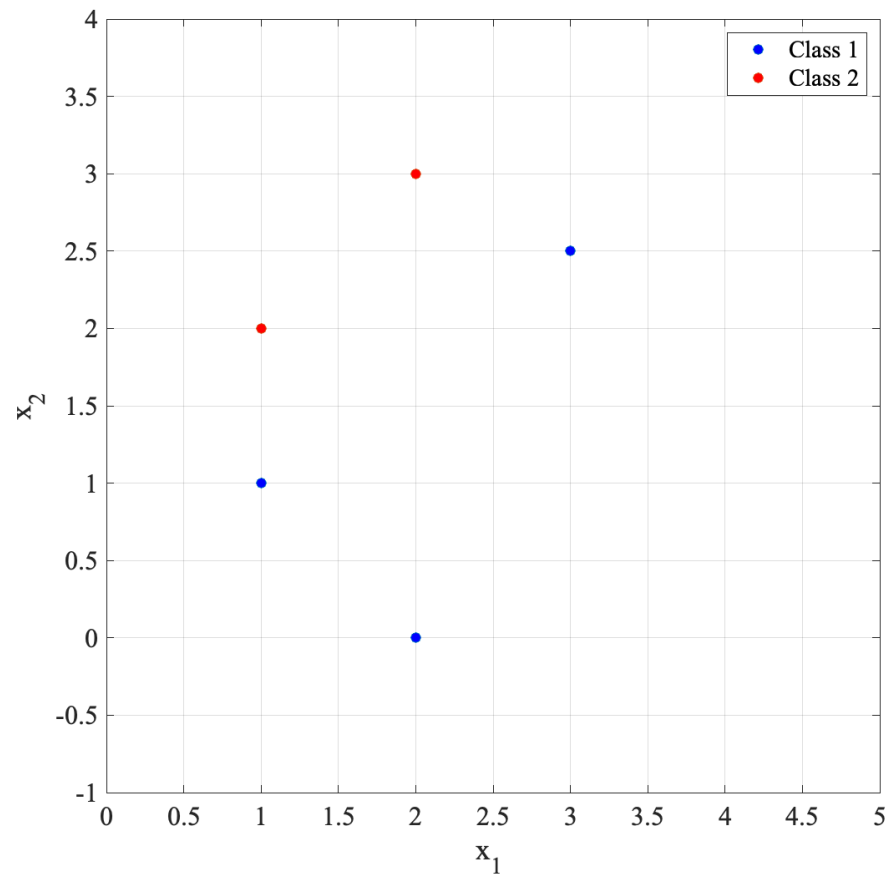
Example

- Class 1:

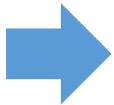
$$\mathbf{x}_1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \mathbf{x}_2 = \begin{bmatrix} 2 \\ 0 \end{bmatrix} \quad \mathbf{x}_3 = \begin{bmatrix} 3 \\ 2.5 \end{bmatrix}$$

- Class 2:

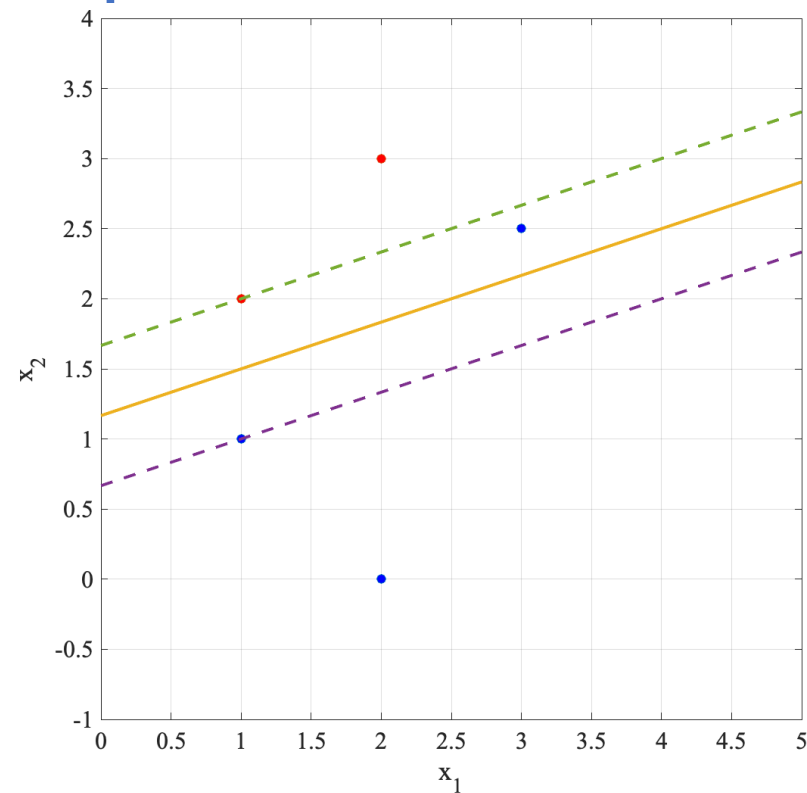
$$\mathbf{x}_4 = \begin{bmatrix} 2 \\ 3 \end{bmatrix} \quad \mathbf{x}_5 = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$



Example

- From the dual problem, using the algorithm for quadratic programming, when $\gamma = 1$, we have $\alpha_1 = 0.16, \alpha_2 = 0.16, \alpha_3 = 0, \alpha_4 = 0$ and $\alpha_5 = 0.32$
- $\begin{bmatrix} 1 \\ 1 \end{bmatrix}, \begin{bmatrix} 2 \\ 0 \end{bmatrix}$ and $\begin{bmatrix} 1 \\ 2 \end{bmatrix}$ are support vectors
- This leads to $w_1 = -0.16$ and $w_2 = 0.48$. This indicates that $w_1 : w_2 = -1 : 3$. We let $\mathbf{w} = [-a, 3a]^T$
- We know that $w^T \begin{bmatrix} 1 \\ 1 \end{bmatrix} + w_0 = -1$ and $w^T \begin{bmatrix} 2 \\ 3 \end{bmatrix} + w_0 = 1$  $a = \frac{2}{3}, w_0 = -1$

The Non-Separable Case



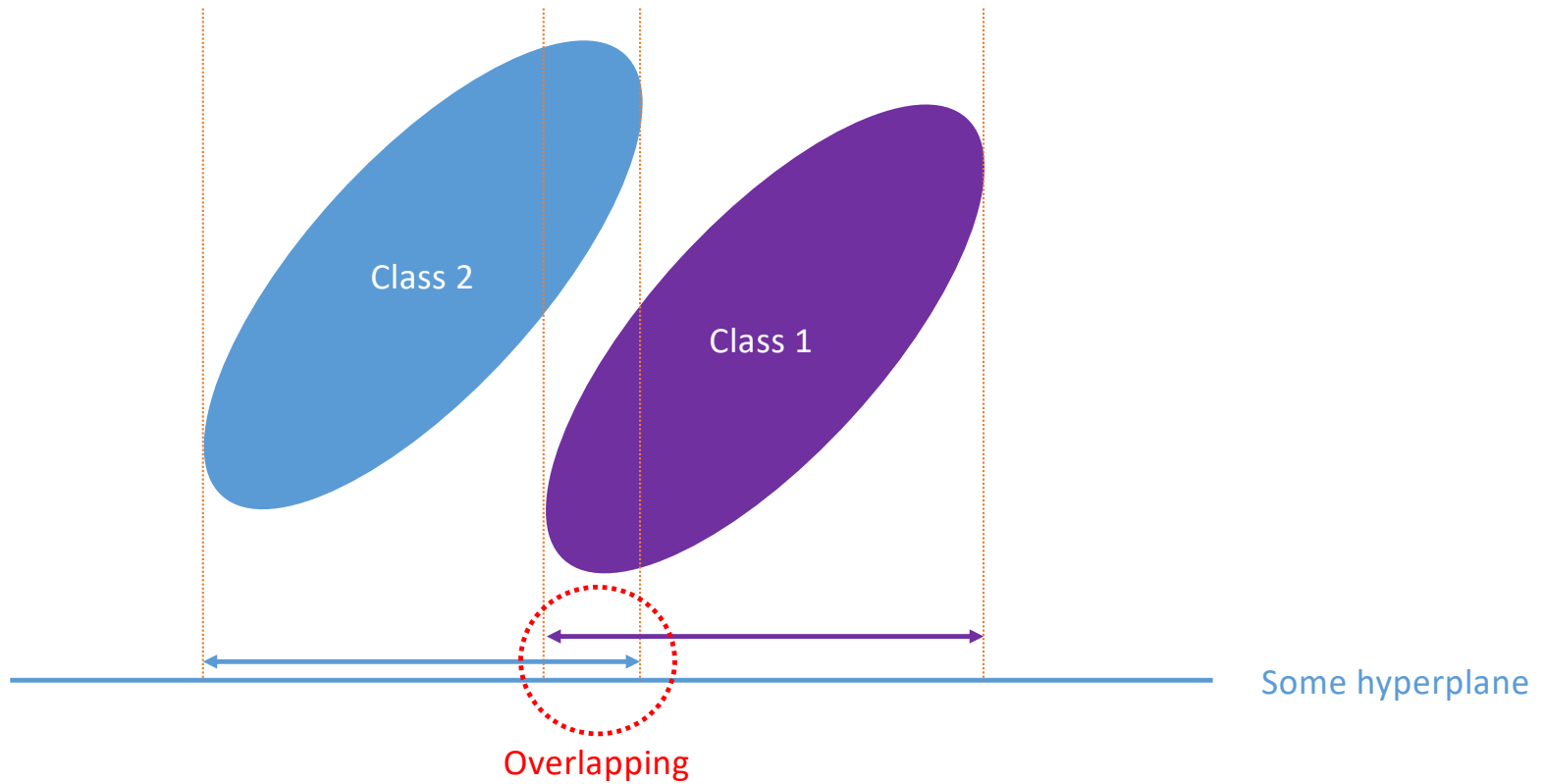
$$\hat{y}\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = -\frac{2}{3}x_1 + \frac{4}{3}x_2 - 1$$

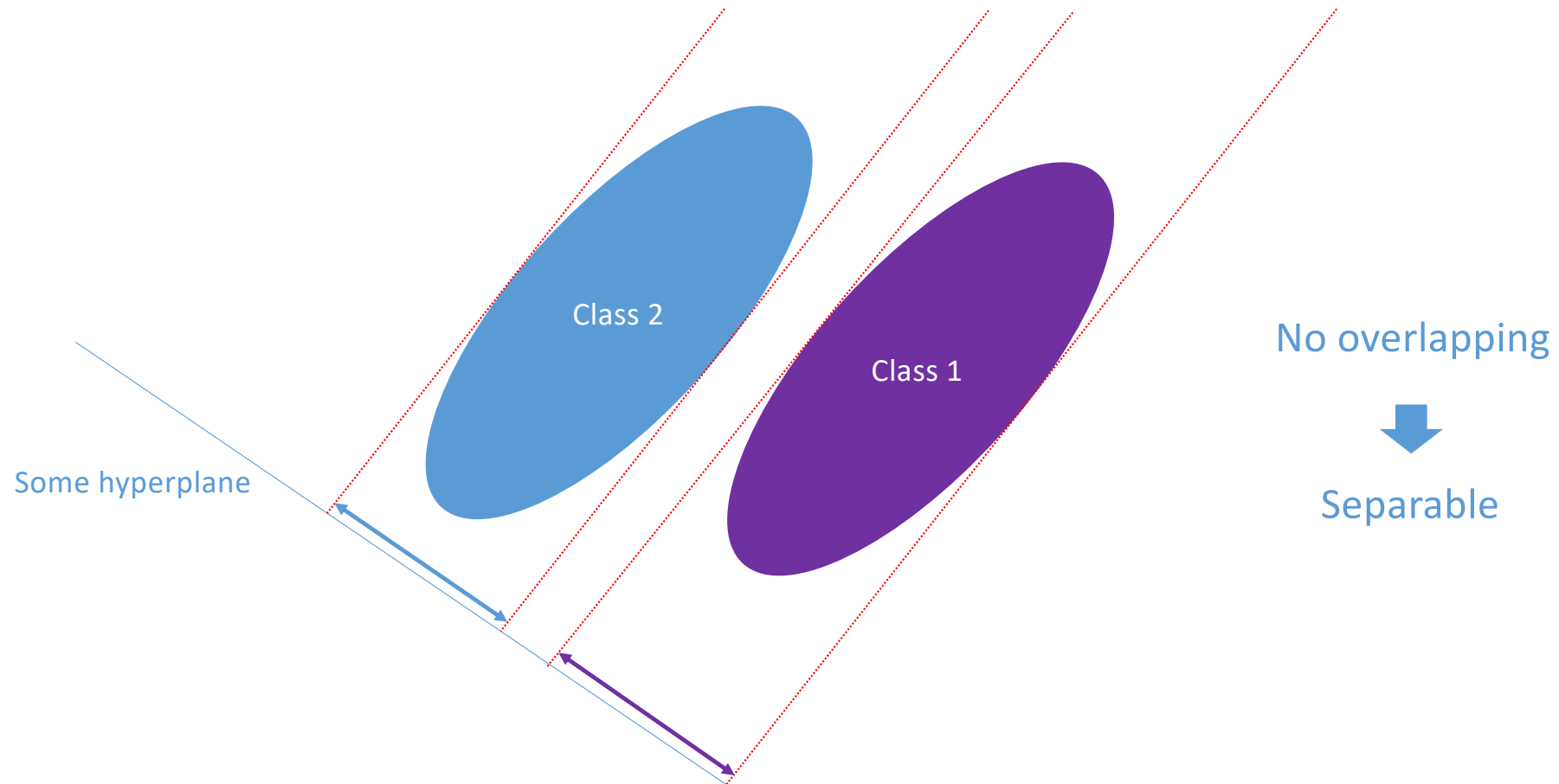
Linear Models for Classification

Topic 05

Fisher's Linear Discriminant

- One way to view a linear classification model is in terms of dimensionality reduction
- The idea is to project the data point onto a hyperplane so that points of different classes can be separated as possible



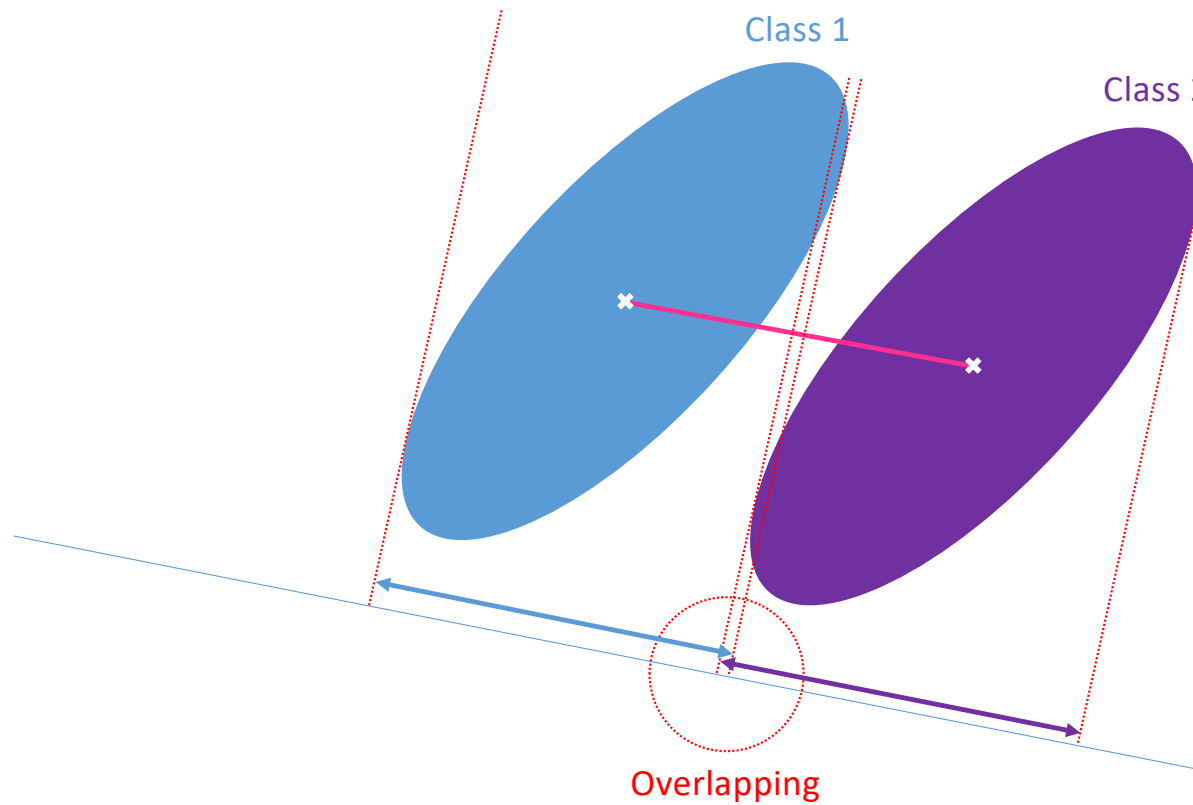


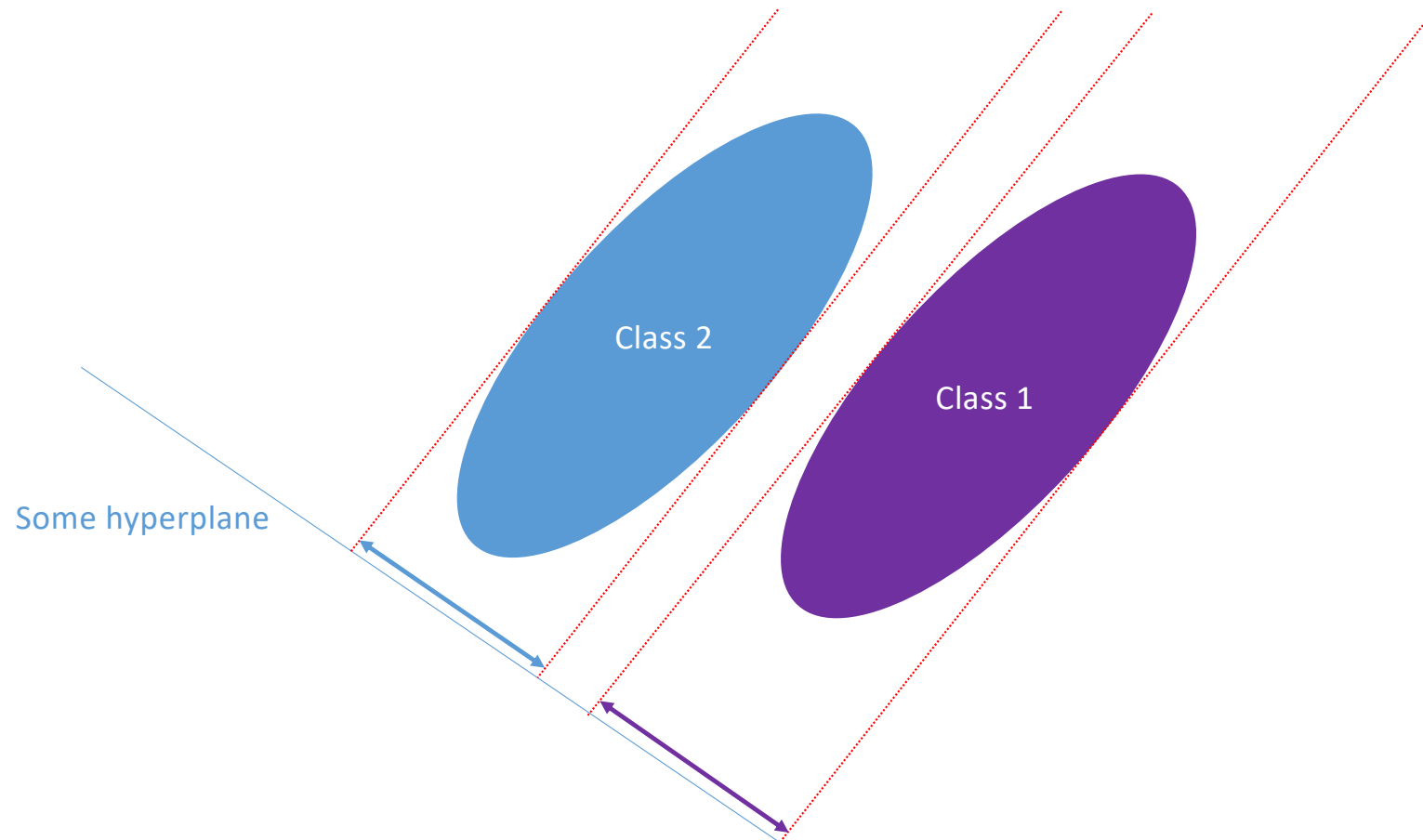
- The goal is to find a projecting vector \mathbf{w} such that these points of these two classes to maximally separate these two classes
- The simplest measure of separation is the the separation of the means of two classes after the projection
- Thus, the goal can be

$$\mathbf{w} = \arg \max_{\mathbf{w}, \|\mathbf{w}\|=1} m_2^p(\mathbf{w}) - m_1^p(\mathbf{w})$$

where $m_k^p(w) = \mathbf{w}^T \mathbf{m}_k$ and \mathbf{m}_k is the mean vector of class C_k for $k = 1, 2$

$$\mathbf{m}_k = \frac{1}{N_k} \sum_{\mathbf{x}_i \in C_k} \mathbf{x}_i$$





字幕預留位置

- The idea of Fisher's Linear Discriminant is to maximize the separation as possible while keeping the dispersion of data points after projection as small as possible
- The goal of Fisher's Linear Discriminant

$$\max_{\mathbf{w}} J(\mathbf{w}) = \frac{\left(m_2^p(\mathbf{w}) - m_1^p(\mathbf{w})\right)^2}{s_1^2(\mathbf{w}) + s_2^2(\mathbf{w})} = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_w \mathbf{w}}$$

Rayleigh quotient

The between-class
covariance matrix

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

The total within-class
covariance matrix

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

- $J(\mathbf{w})$ is maximized when

$$(\mathbf{w}^T \mathbf{S}_B \mathbf{w}) \mathbf{S}_w \mathbf{w} = (\mathbf{w}^T \mathbf{S}_w \mathbf{w}) \mathbf{S}_B \mathbf{w} \rightarrow J(\mathbf{w}) \mathbf{S}_w \mathbf{w} = \boxed{\mathbf{S}_B \mathbf{w}} \rightarrow \lambda \mathbf{w} = \mathbf{S}_w^{-1} \mathbf{S}_B \mathbf{w}$$

the generalized eigenvalue problem

always in the direction of $(\mathbf{m}_2 - \mathbf{m}_1)$

- Let us look at $\mathbf{S}_B \mathbf{w}$ carefully

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

- $\mathbf{S}_B \mathbf{w}$ is always in the same direction of $(\mathbf{m}_2 - \mathbf{m}_1)$
- Since both $\mathbf{w}^T \mathbf{S}_B \mathbf{w}$ and $\mathbf{w}^T \mathbf{S}_w \mathbf{w}$ are scalars, we have

$$\mathbf{S}_w \mathbf{w} \propto \mathbf{S}_B \mathbf{w} = \alpha(\mathbf{m}_2 - \mathbf{m}_1) \Rightarrow \mathbf{w} \propto \boxed{\mathbf{S}_w^{-1}(\mathbf{m}_2 - \mathbf{m}_1)}$$

Fisher's linear discriminant

- After the projecting vector \mathbf{w} is determined, the decision rule can be

$$\hat{y}(\mathbf{x}) = \begin{matrix} C_1 \\ \mathbf{w}^T \mathbf{x} \underset{C_2}{\gtrless} \gamma_{th} \end{matrix}$$

Example – Gaussian data

- Class 1:

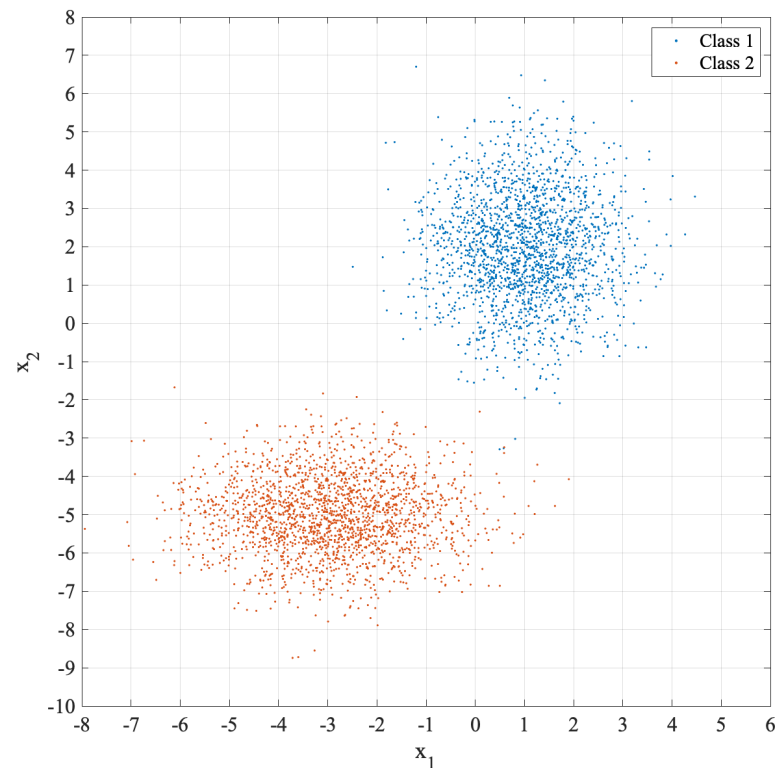
$$\mathbf{m}_1 = \begin{bmatrix} 1 \\ 2 \end{bmatrix} \quad \mathbf{K}_1 = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}$$

- Class 2:

$$\mathbf{m}_2 = \begin{bmatrix} -3 \\ -5 \end{bmatrix} \quad \mathbf{K}_2 = \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix}$$

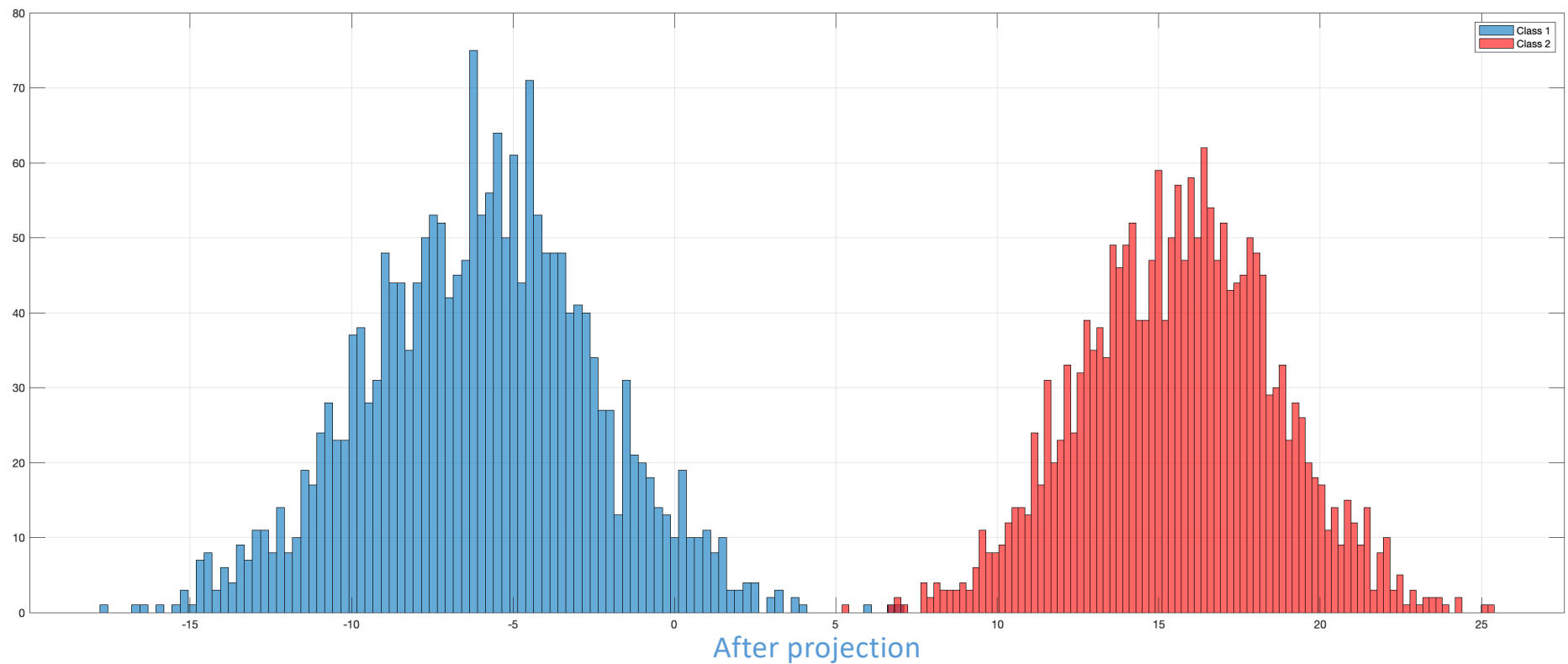
- The projecting vector

$$\begin{aligned} \mathbf{w} &= \mathbf{S}_w^{-1}(\mathbf{m}_2 - \mathbf{m}_1) = (\mathbf{K}_2 + \mathbf{K}_1)^{-1}(\mathbf{m}_2 - \mathbf{m}_1) \\ &= \begin{bmatrix} -4/3 \\ -7/3 \end{bmatrix} \end{aligned}$$





Example – Gaussian data

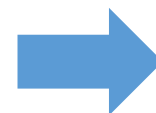


字幕預留位置

Fisher's Discriminant for Multi-class Classification

- Assume that K classes are present and we would like to classify the targets into one of these classes in D -dimensional space, where $D > K$
- Now we introduce D' "features"

$$\hat{y}_i(\mathbf{x}) = \mathbf{w}_i^T \mathbf{x}, \quad \text{for } i = 1, 2, \dots, D'$$



$$\begin{aligned} \hat{\mathbf{Y}} &= \mathbf{W}^T \mathbf{x} \\ \mathbf{W} &= [\mathbf{w}_1 \ \mathbf{w}_2 \ \cdots \ \mathbf{w}_{D'}] \end{aligned}$$

Dimension Reduction

Fisher's Discriminant for Multi-class Classification

- The within-class covariance matrix

$$\mathbf{S}_W = \sum_{k=1}^K \mathbf{S}_k \quad \mathbf{S}_k = \sum_{n \in C_k} (\mathbf{x}_n - \mathbf{m}_k)(\mathbf{x}_n - \mathbf{m}_k)^T$$
$$\mathbf{m}_k = \frac{1}{N_k} \sum_{n \in C_k} \mathbf{x}_n$$

- The between-class covariance

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

- The total covariance matrix

$$\mathbf{S}_T = \mathbf{S}_W + \mathbf{S}_B$$

Fisher's Discriminant for Multi-class Classification

- The goal is similar to the two-class classification is to seek an optimal W to maximize

$$J(W) = \text{tr}\{(W^T S_W W)^{-1} (W^T S_B W)\}$$

- The D' projecting vectors are those eigenvectors corresponding to the D' largest eigenvalues of

$$S_W^{-1} S_B$$

Q&A