

# 第五章 线性学习机器与线性分类器

苏智勇

可视计算研究组

南京理工大学

[suzhiyong@njust.edu.cn](mailto:suzhiyong@njust.edu.cn)

<https://zhiyongsu.github.io>

# 主要内容

5.1 引言

5.2 线性回归

5.3 线性判别函数的基本概念

5.4 Fisher线性判别分析

5.5 感知器

5.6 最小平方误差判别

5.7 罗杰斯特回归

5.8 最优分类超平面与线性支持向量机

5.9 多类线性分类器

# 5.1 引言

- 模式分类的三种主要途径
  - 估计类条件概率密度  $P(x | \omega_i)$ 
    - 通过  $P(\omega_i)$  和  $P(x | \omega_i)$ ，利用贝叶斯规则计算后验概率  $P(\omega_i | x)$ ，然后通过最大后验概率做出决策。
    - 概率密度参数估计和非参数估计。
  - 直接估计后验概率  $P(\omega_i | x)$ 

不需要先估计类条件概率密度  $P(x | \omega_i)$ 。主要有K近邻方法等。
  - **直接计算判别函数**
    - 不需要估计  $P(x | \omega_i)$  或者  $P(\omega_i | x)$ 。
    - 直接找到可用于分类的判别函数。常见的方法有神经网络等。

# 5.1 引言

- **基于概率密度（估计）的分类器设计：model-based method**
  - **已知：** 类先验概率 $P(\omega_i)$ 和类条件概率密度函数 $P(x | \omega_i)$
  - **任务：** 估计一个决策函数，借此进行分类
  - **方法：** 参数估计、非参数估计
  - **特点：** 需要大量的样本，需要知道某些概率及其形式

# 5.1 引言

- 基于样本的直接分类器设计: **data-driven method**

- 思路:

- 如果知道判别函数的形式, 则可以直接从样本估计函数参数
- 先选定判别函数类和一定的目标 (准则), 利用样本集确定函数类中的未定参数, 使所选的准则最好

- 形式化:

- 判别函数类:  $\{g(\alpha), \alpha \in \Lambda\}$ ,  $\alpha$ 为未定参数
- 准则函数/目标函数:  $L(\alpha)$
- 求 $\alpha^*$ : 
$$L(\alpha^*) = \min_{\alpha} L(\alpha)$$

# 5.1 引言

- **本章利用样本直接设计分类器的基本思想**
  - 给定一个判别函数，且已知该函数的参数形式
  - 采用样本来训练判别函数的参数
  - 对于新样本，采用判别函数对其进行判决，并按照一些准则来完成分类
- **本章利用样本直接设计分类器的基本技术路线**
  - 假定有  $n$  个  $d$  维空间中的样本，每个样本的类别标签已知，且一共有  $c$  个不同的类别。
  - 假定判别函数的形式已知，寻找一个判别函数。
  - 对于给定的新样本  $\mathbf{x} \in \mathbb{R}^d$ ，判定它属于  $\omega_1, \omega_2, \dots, \omega_c$  中的哪个类别。

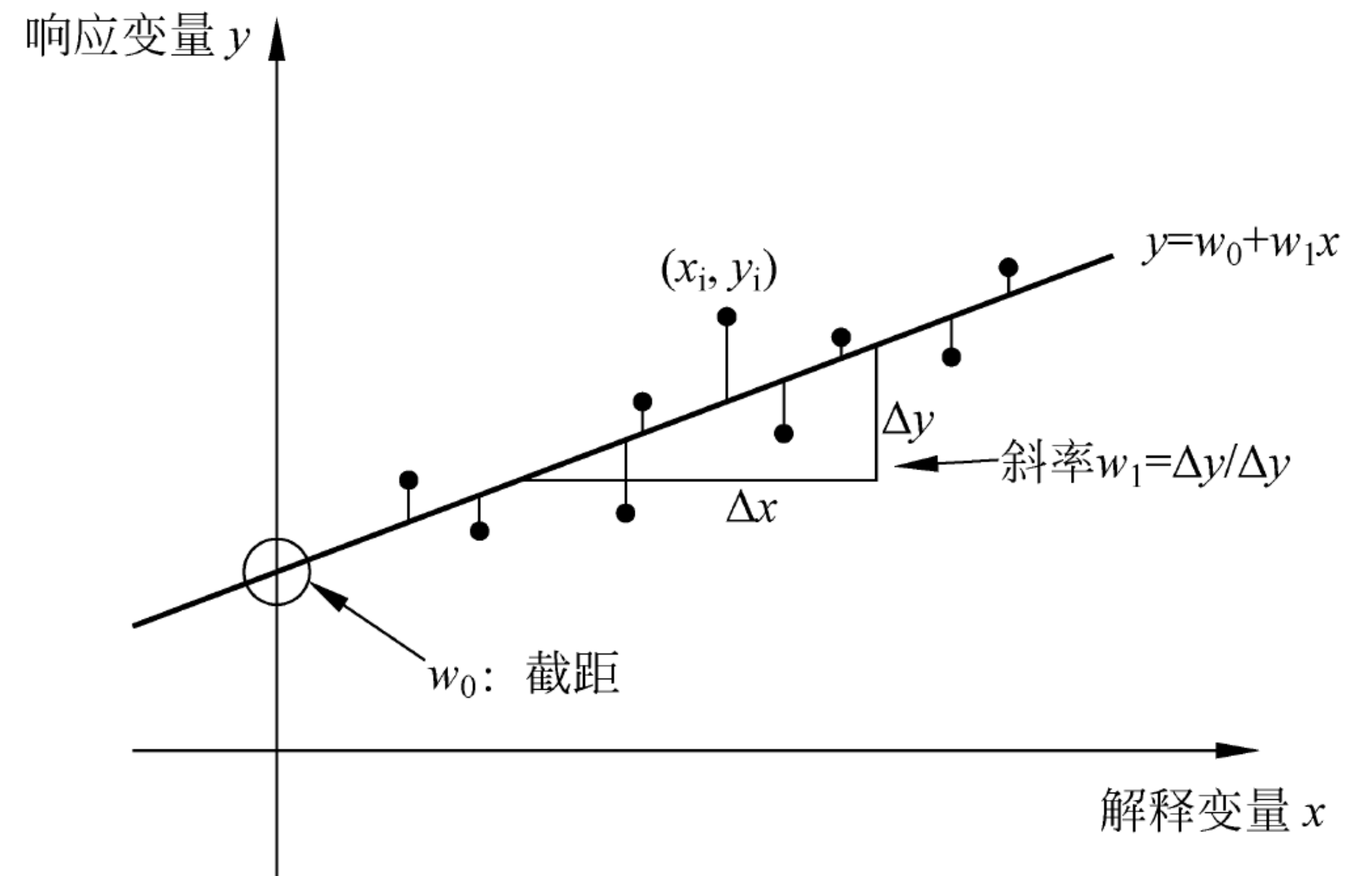
# 5.2 线性回归

- 线性回归

- 通过数据发现或估计两个或多个变量之间可能存在的线性依赖关系的基本统计学方法

- 简单线性回归

- 解释变量  $x$ 、响应变量  $y$
- 通过  $(x, y)$  的一系列观测样本，估计线性关系  $y = w_0 + w_1x$ ，即估计其中的系数  $w_0$  和  $w_1$



# 5.2 线性回归

- 多元线性回归

— 响应变量依赖于多个解释变量:  $y = w_0 + w_1x + \dots + w_dx_d = \sum_{i=0}^d w_ix_i = \mathbf{w}^T \mathbf{x}$

- 线性回归问题描述

— 已知训练样本集  $\left\{ (x_1, y_1), \dots, (x_N, y_N) \right\}$ ,  $x_j \in R^{d+1}$ ,  $y_j \in R$ ; 机器学习模型为

$$f(\mathbf{x}) = w_0 + w_1x + \dots + w_dx_d = \sum_{i=0}^d w_ix_i = \mathbf{w}^T \mathbf{x}, \text{ 其中,}$$

$\mathbf{w} = [w_0, w_1, \dots, w_d]^T$  是模型中的待定参数



# 5.2 线性回归

- “最小二乘法” 求解线性回归问题

- 用训练样本集估计模型中的参数，使模型在最小平方误差意义下能够最好地拟合训练样

- 本，即 
$$\min E = \frac{1}{N} \sum_{j=1}^N (f(x_j) - y_j)^2$$

- 目标函数可进一步写成矩阵形式 
$$E(\mathbf{w}) = \frac{1}{N} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 = \frac{1}{N} (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}),$$

其中， $\mathbf{X} = \begin{bmatrix} x_1^T \\ \vdots \\ x_N^T \end{bmatrix}$  是全部训练样本的解释变量向量组成的矩阵， $\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}$  是全部训练样

本的响应变量组成的向量

# 5.2 线性回归

- “最小二乘法” 求解线性回归问题

- 使目标函数 $E(\mathbf{w})$ 最小化的参数 $\mathbf{w}$ 应满足 $\frac{\partial E(\mathbf{w})}{\partial(\mathbf{w})} = \frac{2}{N} \mathbf{X}^T (\mathbf{X}\mathbf{w} - \mathbf{y}) = 0$ ,

- 即  $\mathbf{X}^T \mathbf{X}\mathbf{w} = \mathbf{X}^T \mathbf{y}$

- 当矩阵 $(\mathbf{X}^T \mathbf{X})$ 可逆时, 最优参数的解为 $\mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ ,

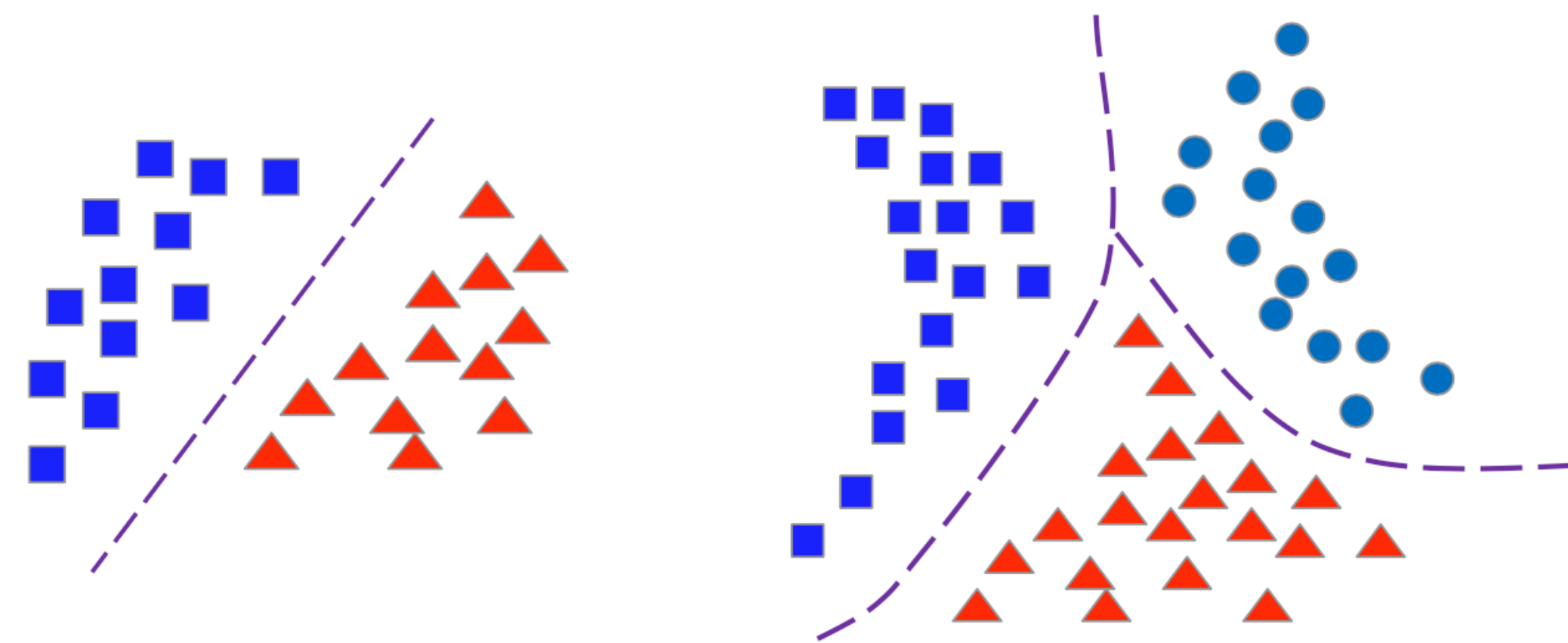
- $(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ 被称为 $\mathbf{X}$ 的伪逆 (pseudo-inverse) 矩阵, 记作 $\mathbf{X}^+$

- 线性回归给出了在最小平方误差意义下对解释变量与响应变量间线性关系的最好估计

# 5.3 线性判别函数的基本概念

- 基于判别函数的分类器

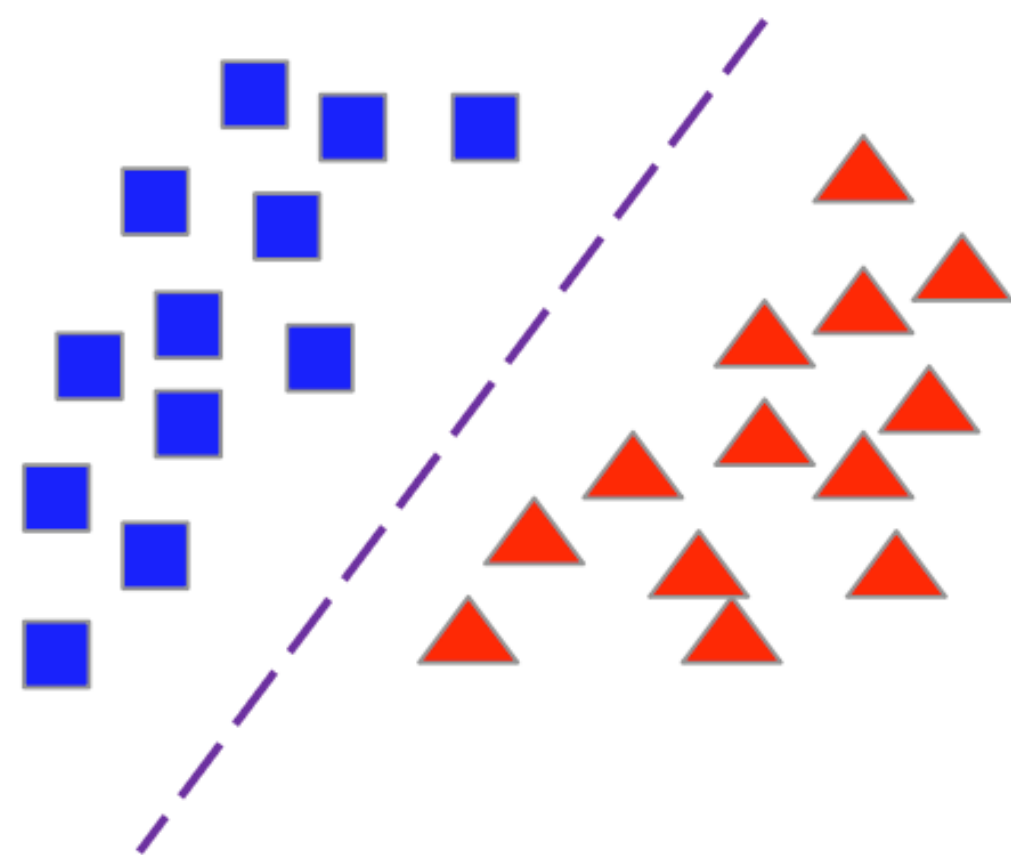
- 采用已知类别标签的训练样本进行学习，获得若干个代数界面，这些界面将样本所在的空间分成若干个相互不重叠的区域。每个区域包含属于同一类的样本。
- 表示界面的函数称为判别函数。
- 判别函数是分类器最常用的表述形式。



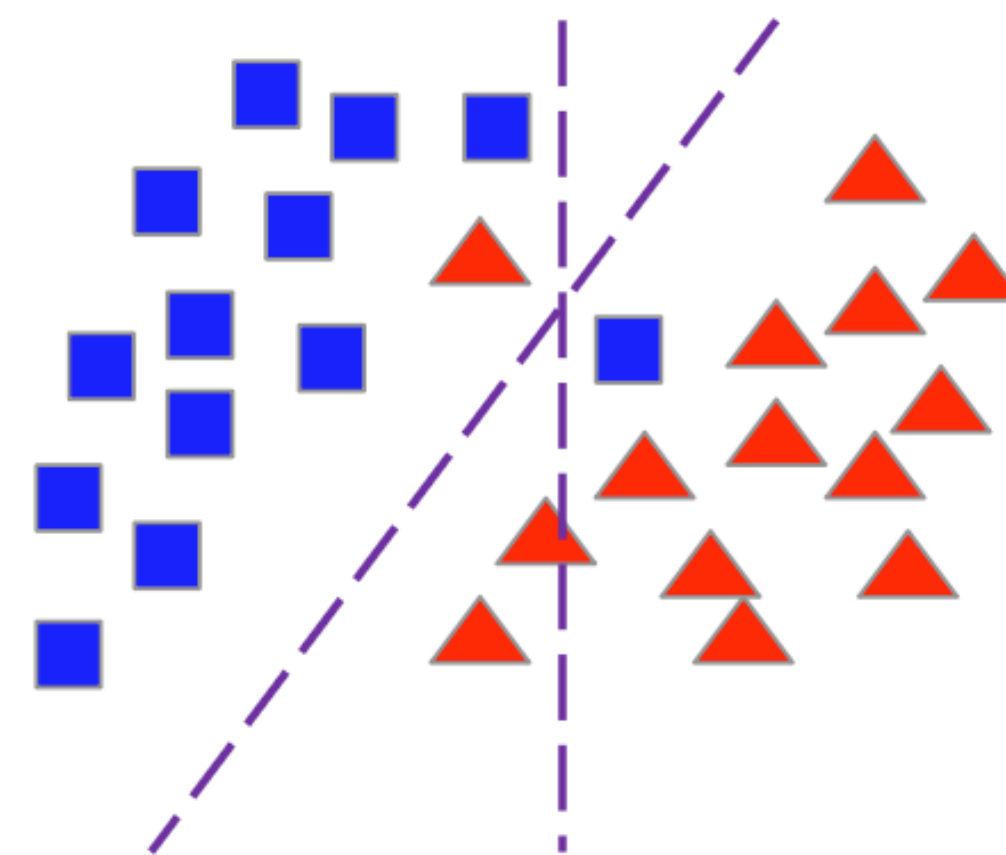
# 5.3 线性判别函数的基本概念

- 线性可分

- 对于  $n$  个  $d$  维空间中的样本  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ , 假定这些样本来自两个类别  $\omega_1$  和  $\omega_2$ 。如果存在一个**线性**判别函数能对这些样本正确地分类, 则称这些样本是**线性可分**的; 否则是**线性不可分**的。



线性可分



线性不可分

# 5.3 线性判别函数的基本概念

- 一般表达式

- $g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$

- $\mathbf{x}$ :  $d$ 维特征向量 (样本向量),  $\mathbf{x} = [x_1, x_2, \dots, x_d]^T$

- $\mathbf{w}$ : 权向量,  $\mathbf{w} = [w_1, w_1, \dots, w_d]^T$

- 对于 $c$ 类分类问题,  $g_i(\mathbf{x})$ ,  $i = 1, 2, \dots, c$ , 表示每个类别对应的判别函数

- 决策规则: 如果 $g_i(\mathbf{x}) > g_j(\mathbf{x})$ ,  $\forall j \neq i$ , 则 $\mathbf{x} \in \omega_i$

# 5.3 线性判别函数的基本概念

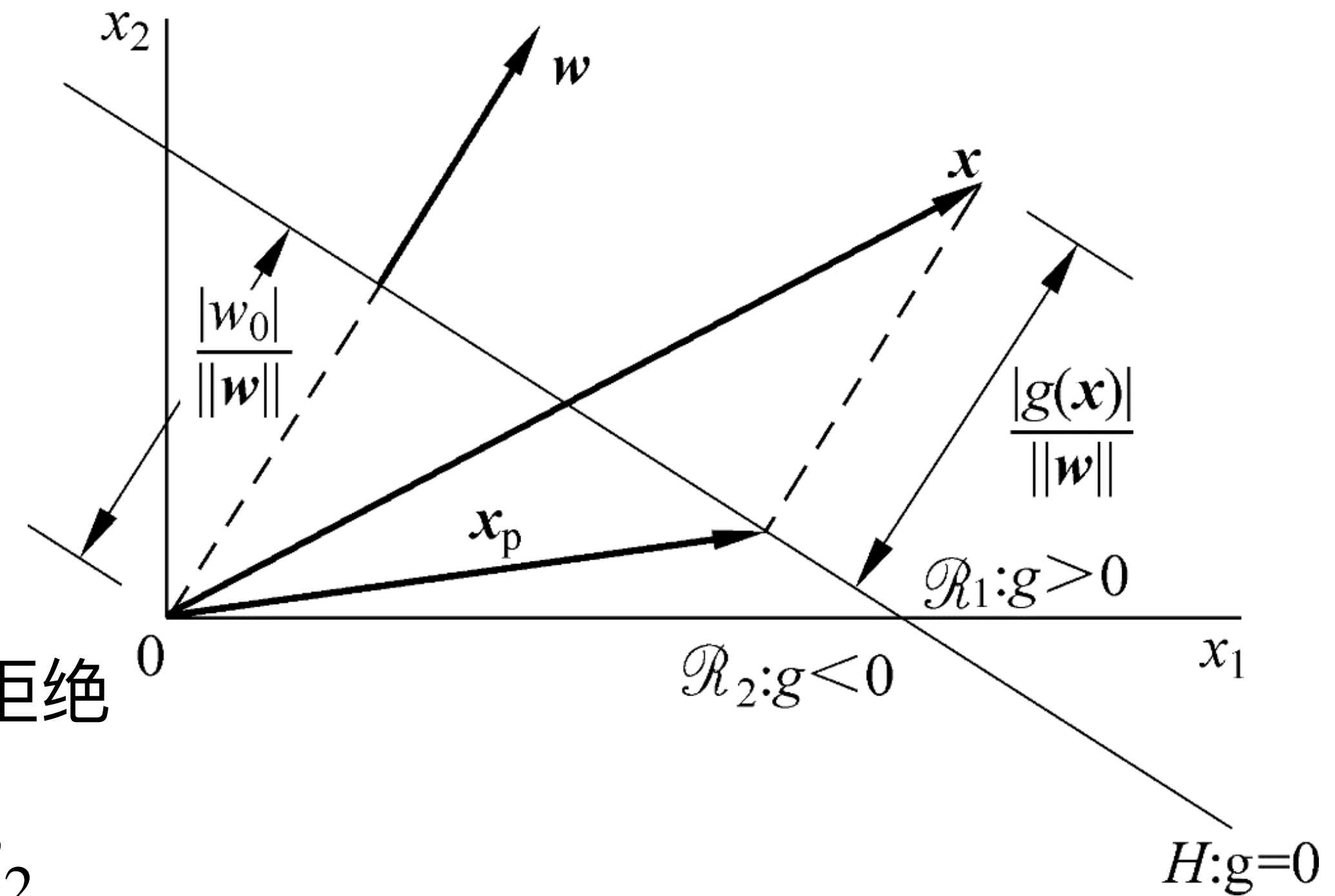
## • 决策规则

– 以二分类问题为例，令  $g(\mathbf{x}) = g_1(\mathbf{x}) - g_2(\mathbf{x})$

如果  $g(\mathbf{x}) > 0$ ，则决策  $\mathbf{x} \in \omega_1$

如果  $g(\mathbf{x}) < 0$ ，则决策  $\mathbf{x} \in \omega_2$

如果  $g(\mathbf{x}) = 0$ ，可将  $\mathbf{x}$  任意分到某一类，或拒绝



✓ 决策面/超平面  $H: g(\mathbf{x}) = 0$ ，决策区域  $\mathcal{R}_1$  和  $\mathcal{R}_2$

✓  $\mathbf{w}$  为超平面  $H$  的法向量

▸ 设  $\mathbf{x}_1, \mathbf{x}_2$  都在  $H$  上，则  $\mathbf{w}^T \mathbf{x}_1 + w_0 = \mathbf{w}^T \mathbf{x}_2 + w_0$ ，即  $\mathbf{w}^T (\mathbf{x}_1 - \mathbf{x}_2) = 0$

▸  $\mathbf{w}^T$  跟决策面  $H$  上的任一向量正交，即  $\mathbf{w}$  为超平面  $H$  的法向量

# 5.3 线性判别函数的基本概念

- 代数度量

- 判别函数 $g(\mathbf{x})$ 正比于 $\mathbf{x}$ 点到超平面 $H$ 的代数距离（带正负号）， $\mathbf{x}$ 在 $H$ 正侧：

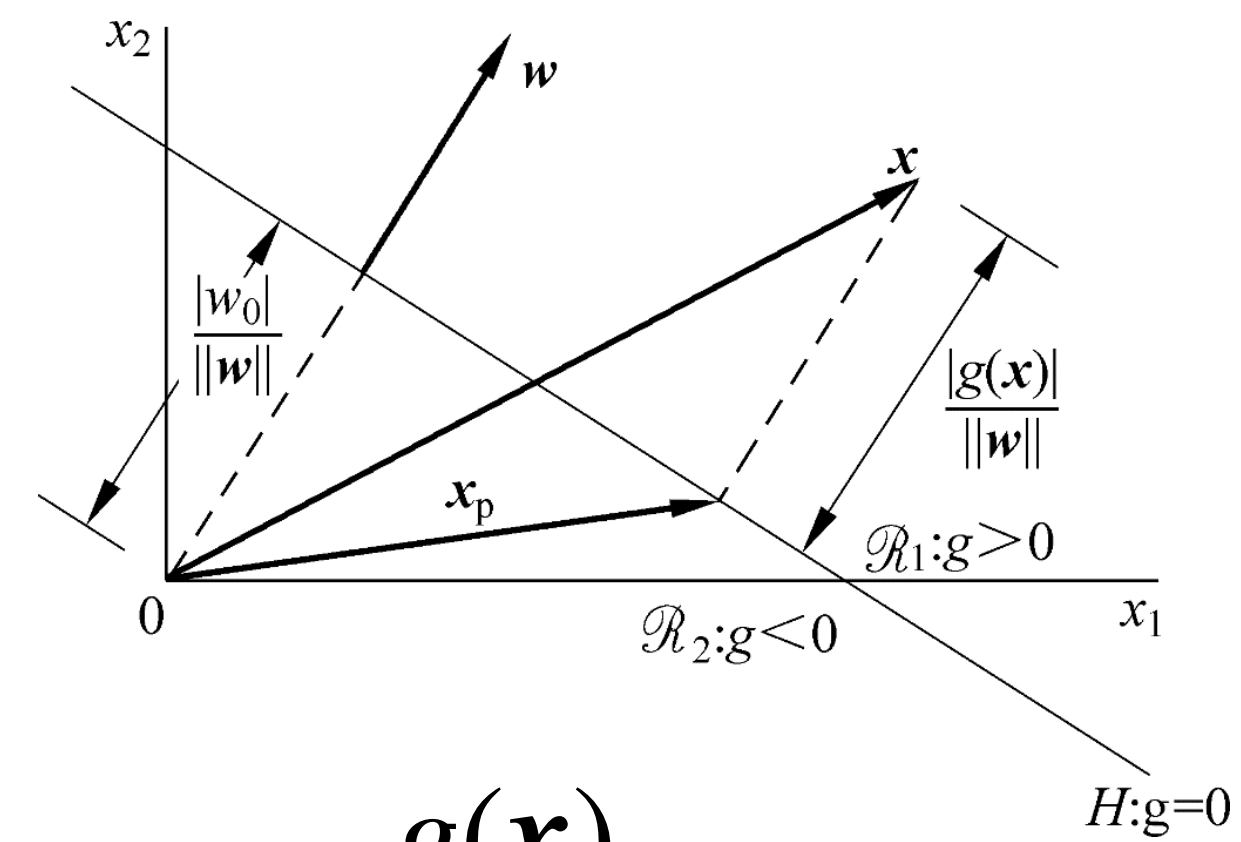
- $g(\mathbf{x}) > 0$ ；在负侧： $g(\mathbf{x}) < 0$

- $g(\mathbf{x})$ 可以视为特征空间中某点 $\mathbf{x}$ 到超平面 $H$ 的距离的一种代数度量

- $\mathbf{x} = \mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|}$ ， $\mathbf{x}_p$ 为 $\mathbf{x}$ 在 $H$ 上的投影向量， $\frac{\mathbf{w}}{\|\mathbf{w}\|}$ 为单位向量

- $g(\mathbf{x}) = \mathbf{w}^\top \left( \mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|} \right) + w_0 = \mathbf{w}^\top \mathbf{x}_p + w_0 + r \frac{\mathbf{w}^\top \mathbf{w}}{\|\mathbf{w}\|} = r \|\mathbf{w}\|$  或  $r = \frac{g(\mathbf{x})}{\|\mathbf{w}\|}$

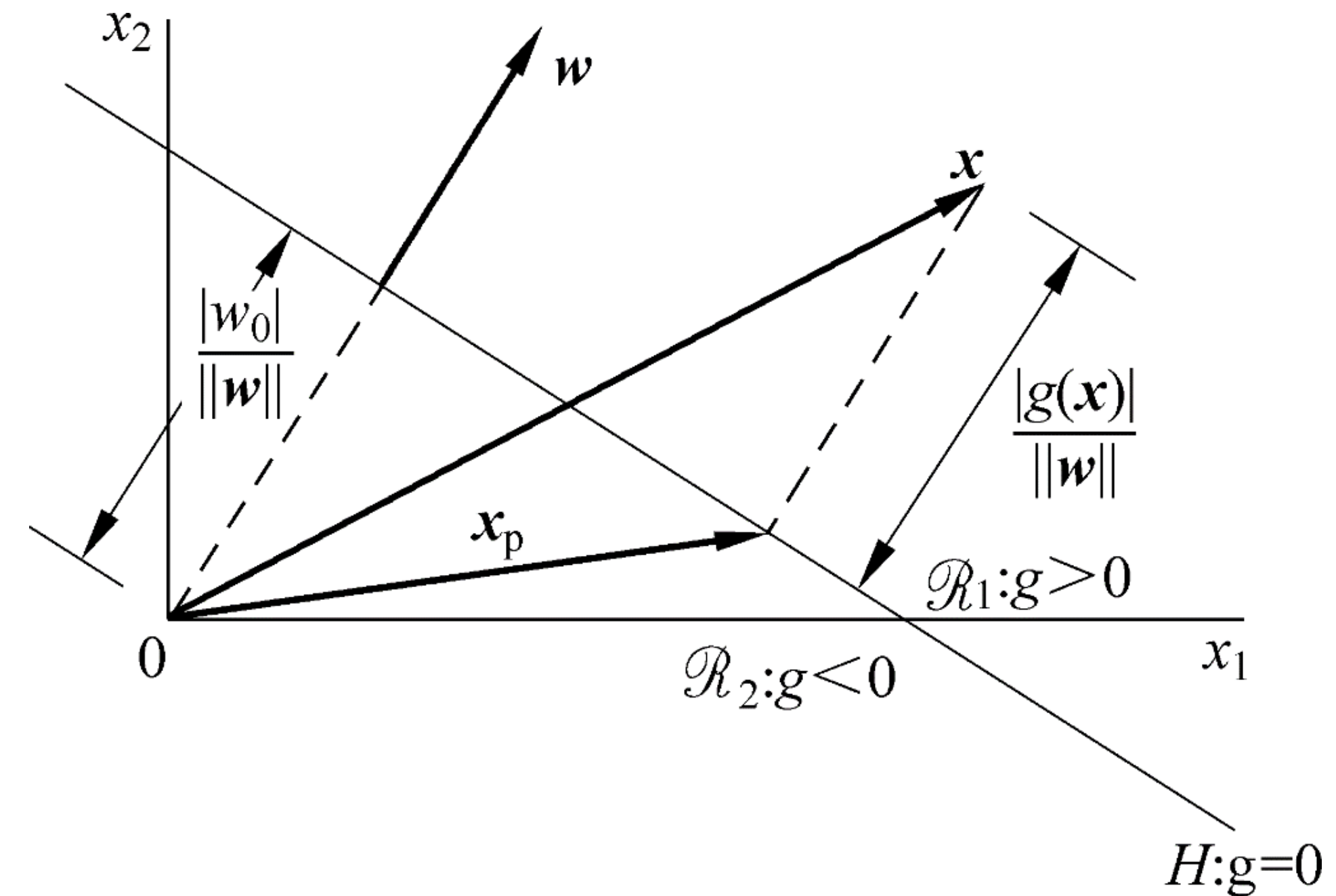
- 若 $\mathbf{x}$ 为原点，则 $g(\mathbf{x}) = w_0$ ， $r_0 = \frac{w_0}{\|\mathbf{w}\|}$



# 5.3 线性判别函数的基本概念

## • 小结

- 用线性判别函数进行决策，就是用一个超平面  $H$  把特征空间分割成决策区域  $\mathcal{R}_1$  和  $\mathcal{R}_2$
- 超平面  $H$  的方向由权向量  $\mathbf{w}$  确定，它的位置由阈值  $w_0$  确定
- 判别函数  $g(\mathbf{x})$  正比于  $\mathbf{x}$  点到超平面的代数距离（带正负号）。当  $\mathbf{x}$  在  $H$  正侧时， $g(\mathbf{x}) > 0$ ；在负侧时， $g(\mathbf{x}) < 0$ 。

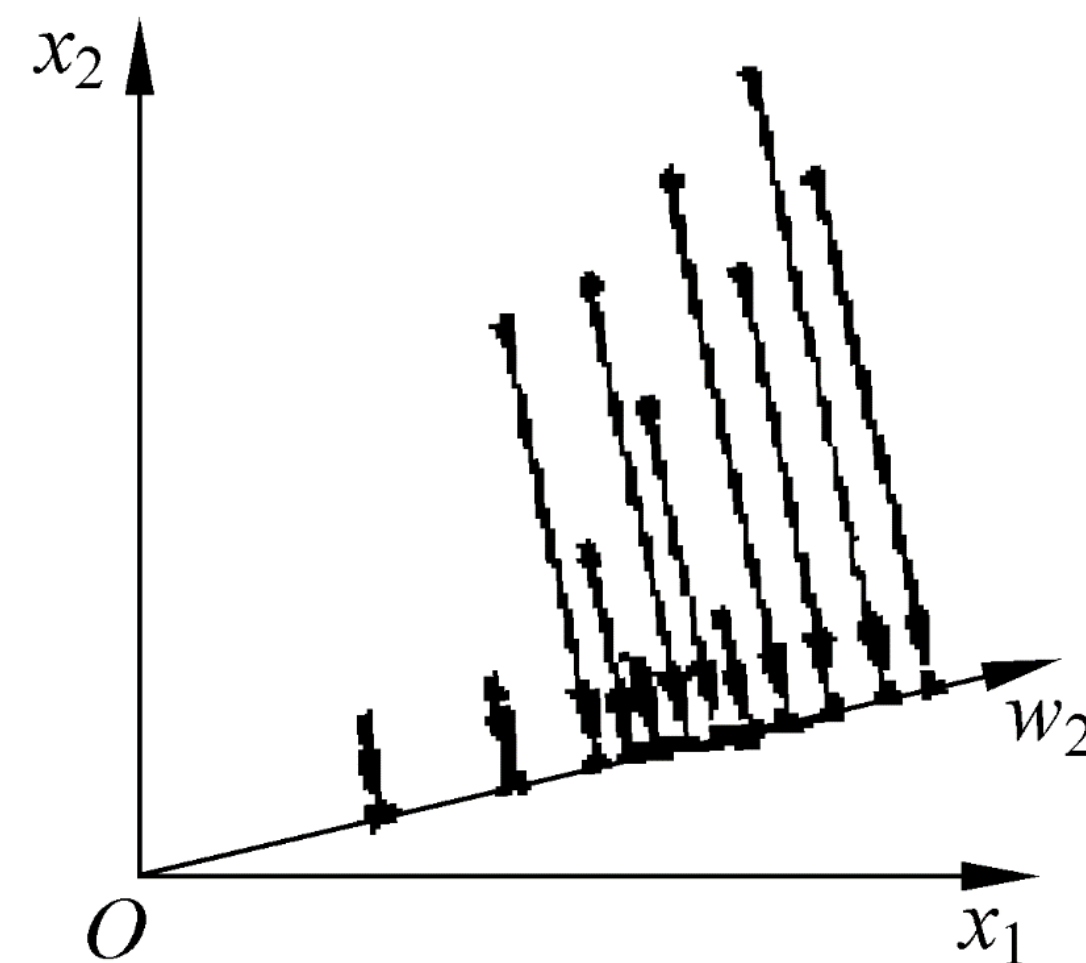
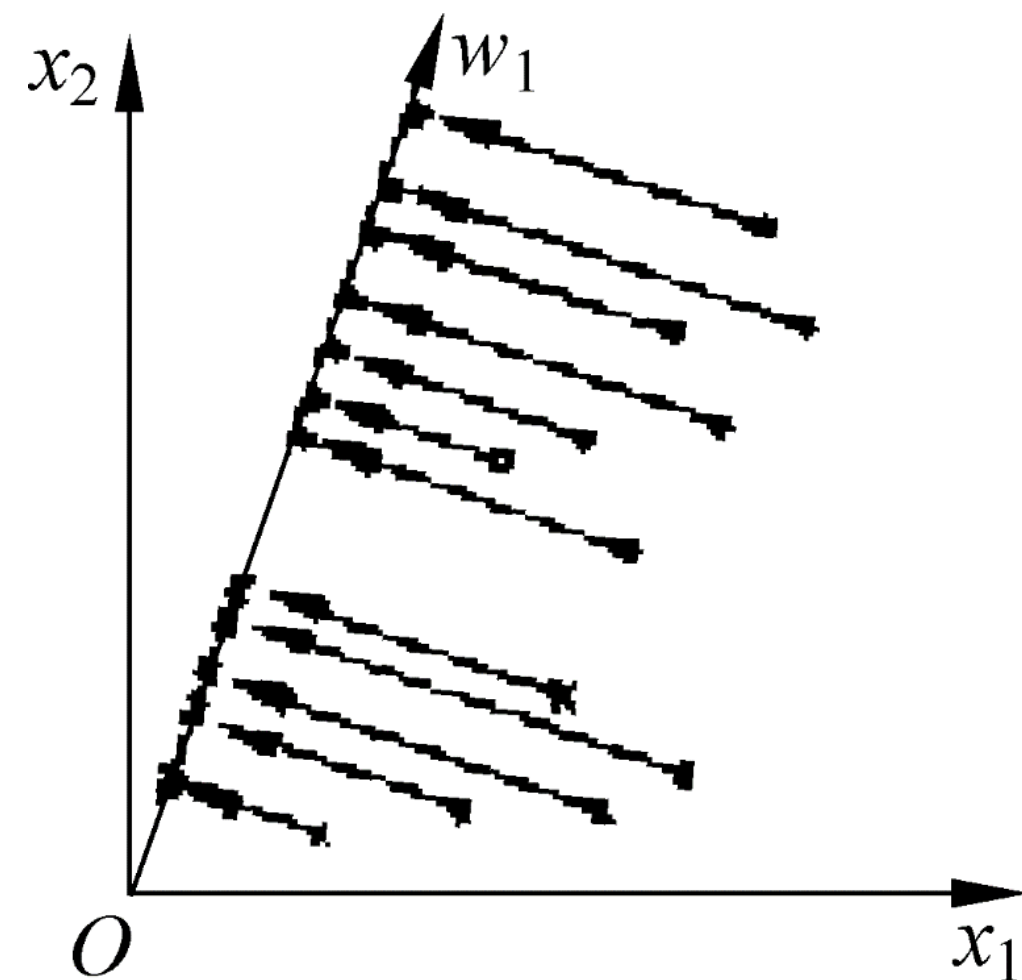




# 5.4 Fisher线性判别分析

- 思路

- Fisher在1936年提出线性判别分析法(Linear Discriminant Analysis, LDA)
- 把所有样本都投影到一个方向上, 在一维空间中确定一个分类的阈值
- 核心: 确定投影方向, 使得类间距离尽可能大, 类内距离尽可能小



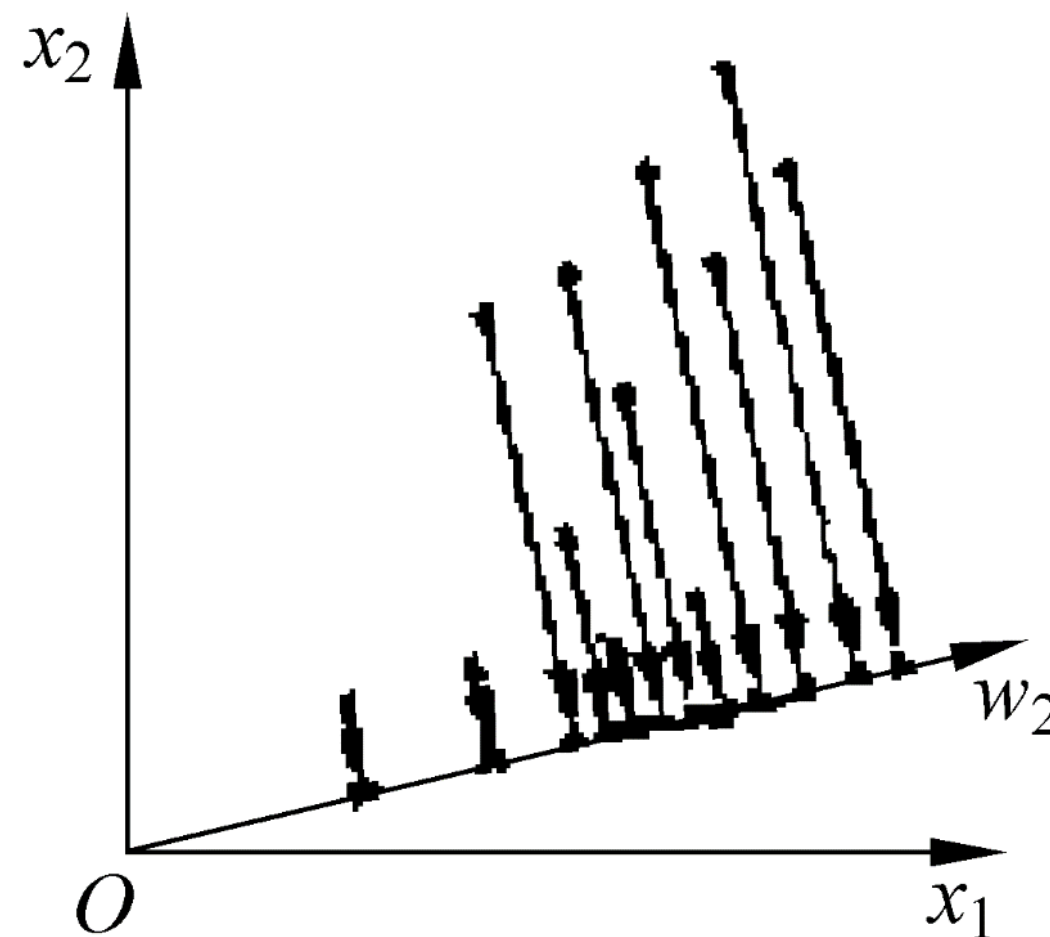
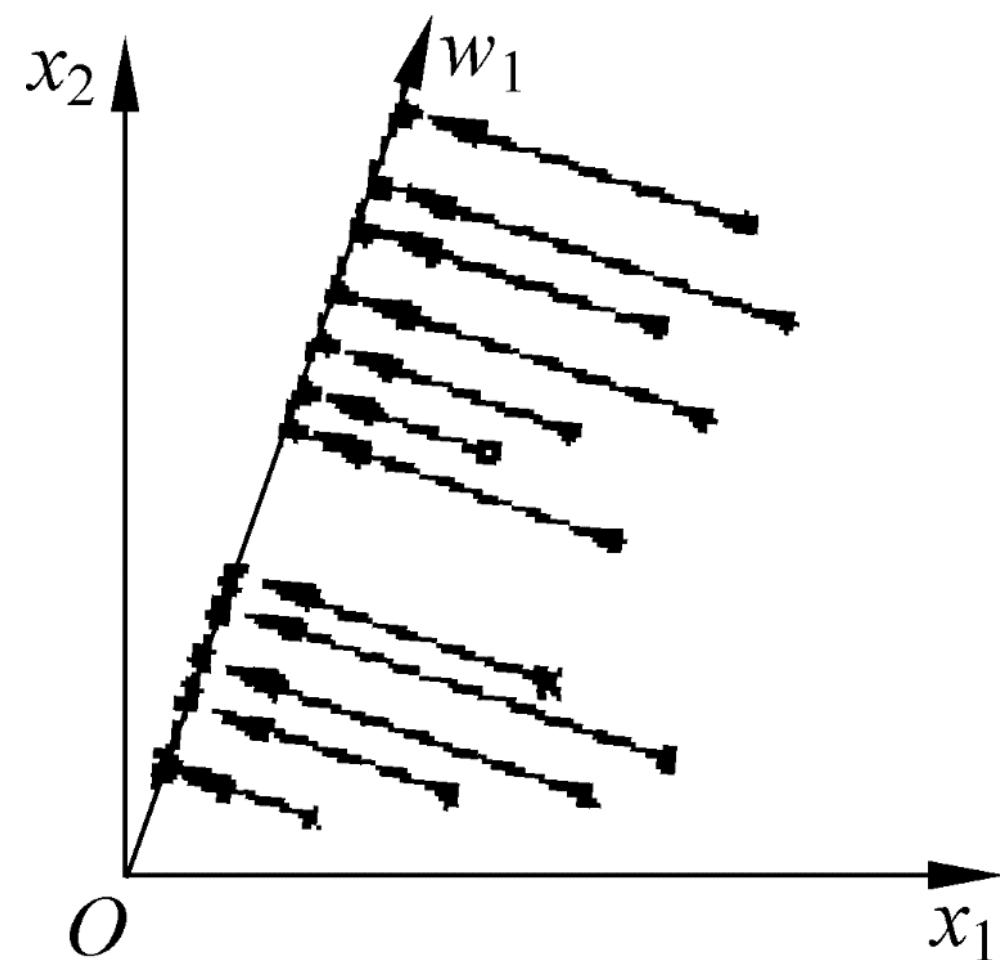
# 5.4 Fisher线性判别分析

- 基本定义

- 以两分类问题为例，训练样本集  $X = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ ，每个样本为  $d$  维向量

- $\omega_1$  类的样本  $X_1 = \{\mathbf{x}_1^1, \dots, \mathbf{x}_{N_1}^1\}$ ， $\omega_2$  类的样本  $X_2 = \{\mathbf{x}_1^2, \dots, \mathbf{x}_{N_2}^2\}$

- 寻找投影方向  $\mathbf{w}$ ，使得投影后的样本为  $y_i = \mathbf{w}^T \mathbf{x}_i, i = 1, 2, \dots, N$



# 5.4 Fisher线性判别分析

- 投影前

- 类内离散度矩阵 (within-class scatter matrix)

$$S_i = \sum_{x_j \in X_i} (\mathbf{x}_j - \mathbf{m}_i)(\mathbf{x}_j - \mathbf{m}_i)^T, \quad i = 1, 2, \quad \text{其中, 类均值向量为 } \mathbf{m}_i = \frac{1}{N_i} \sum_{x_j \in X_i} \mathbf{x}_j, \quad i = 1, 2$$

- 总类内离散度矩阵 (pooled within-class scatter matrix)

$$S_w = S_1 + S_2$$

- 类间离散度矩阵 (between-class scatter matrix)

$$S_b = (\mathbf{m}_1 - \mathbf{m}_2)(\mathbf{m}_1 - \mathbf{m}_2)^T, \quad i = 1, 2$$

# 5.4 Fisher线性判别分析

- 投影后一维空间

- 类内离散度

$$\tilde{S}_i^2 = \sum_{y_j \in Y_i} (y_j - \tilde{m}_i)^2, \quad \text{其中, 两类均值为 } \tilde{m}_i = \frac{1}{N_i} \sum_{y_j \in Y_i} y_j = \frac{1}{N_i} \sum_{\mathbf{x}_j \in X_i} \mathbf{w}^T \mathbf{x}_j = \mathbf{w}^T \mathbf{m}_i, \quad i = 1, 2$$

- 总类内离散度矩阵

$$\tilde{S}_w = \tilde{S}_1^2 + \tilde{S}_2^2$$

- 类间离散度

$$\tilde{S}_b = (\tilde{m}_1 - \tilde{m}_2)^2, \quad i = 1, 2$$

# 5.4 Fisher线性判别分析

- Fisher准则函数 (Fisher's Criterion)

$$\max J_F(\mathbf{w}) = \frac{\tilde{S}_b}{\tilde{S}_w} = \frac{(\tilde{m}_1 - \tilde{m}_2)^2}{\tilde{S}_1^2 + \tilde{S}_2^2}, \text{ 即, 使投影后两类尽可能分开, 而各}$$

类内部尽可能聚集

$$\cdot \tilde{S}_b = (\tilde{m}_1 - \tilde{m}_2)^2 = (\mathbf{w}^T \mathbf{m}_1 - \mathbf{w}^T \mathbf{m}_2)^2 = \mathbf{w}^T \mathbf{S}_b \mathbf{w}$$

$$\cdot \tilde{S}_w = \tilde{S}_1^2 + \tilde{S}_2^2 = \mathbf{w}^T \mathbf{S}_1 \mathbf{w} + \mathbf{w}^T \mathbf{S}_2 \mathbf{w} = \mathbf{w}^T \mathbf{S}_w \mathbf{w}$$

$$\max J_F(\mathbf{w}) = \frac{\mathbf{w}^T \mathbf{S}_b \mathbf{w}}{\mathbf{w}^T \mathbf{S}_w \mathbf{w}}: \text{ 广义的瑞利商 (generalized Rayleigh quotient)}$$

# 5.4 Fisher线性判别分析

- 求解

- 上述 Fisher 准则函数转化为:  $\max \mathbf{w}^T \mathbf{S}_b \mathbf{w} \quad \text{s.t.} \quad \mathbf{w}^T \mathbf{S}_w \mathbf{w} = c \neq 0$

- 定义拉格朗日函数  $L(\mathbf{w}, \lambda) = \mathbf{w}^T \mathbf{S}_b \mathbf{w} - \lambda(\mathbf{w}^T \mathbf{S}_b \mathbf{w} - c)$

- 令  $\frac{\partial L(\mathbf{w}, \lambda)}{\partial \mathbf{w}} = 0$ , 极值解  $\mathbf{w}^*$  应满足:  $\mathbf{S}_b \mathbf{w}^* - \lambda \mathbf{S}_w \mathbf{w}^* = 0$

- 假设  $\mathbf{S}_w$  是非奇异的 ( $N > d$  时通常非奇异), 可得  $\mathbf{S}_w^{-1} \mathbf{S}_b \mathbf{w}^* = \lambda \mathbf{w}^*$ , 即  $\mathbf{w}^*$  为  $\mathbf{S}_w^{-1} \mathbf{S}_b$  的本征向量

- 将  $\mathbf{S}_b$  带入可得  $\lambda \mathbf{w}^* = \mathbf{S}_w^{-1} (\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w}^*$ , 只考虑  $\mathbf{w}^*$  的方向, 得最优投影方向:  $\mathbf{w}^* = \mathbf{S}_w^{-1} (\mathbf{m}_1 - \mathbf{m}_2)$ : 只需要求出原始样本的均值和方差!

# 5.4 Fisher线性判别分析

- 决策规则

- Fisher判别函数的最优解只给出了一个投影方向，不是分类面
- 确定分类阈值 $w_0$ ，得到分类面和决策规则

$$g(x) = \mathbf{w}^T \mathbf{x} + w_0 \lesseqgtr 0, \text{ 则 } x \in \begin{cases} \omega_1 \\ \omega_2 \end{cases}$$

# 5.4 Fisher线性判别分析

- 确定分类阈值 $w_0$

- 样本正态分布且两类协方差矩阵相同时

- 最优贝叶斯分类器为线性函数 $g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$ ,  $\mathbf{w} = \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$ ,

- $$w_0 = -\frac{1}{2} (\boldsymbol{\mu}_1 + \boldsymbol{\mu}_2)^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) - \ln \frac{P(\omega_2)}{P(\omega_1)}$$

- 对于Fisher线性判别:

- $$w_0 = -\frac{1}{2} (\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{S}_w^{-1} (\mathbf{m}_1 - \mathbf{m}_2) - \ln \frac{P(\omega_2)}{P(\omega_1)}$$



# 5.4 Fisher线性判别分析

- 确定分类阈值 $w_0$

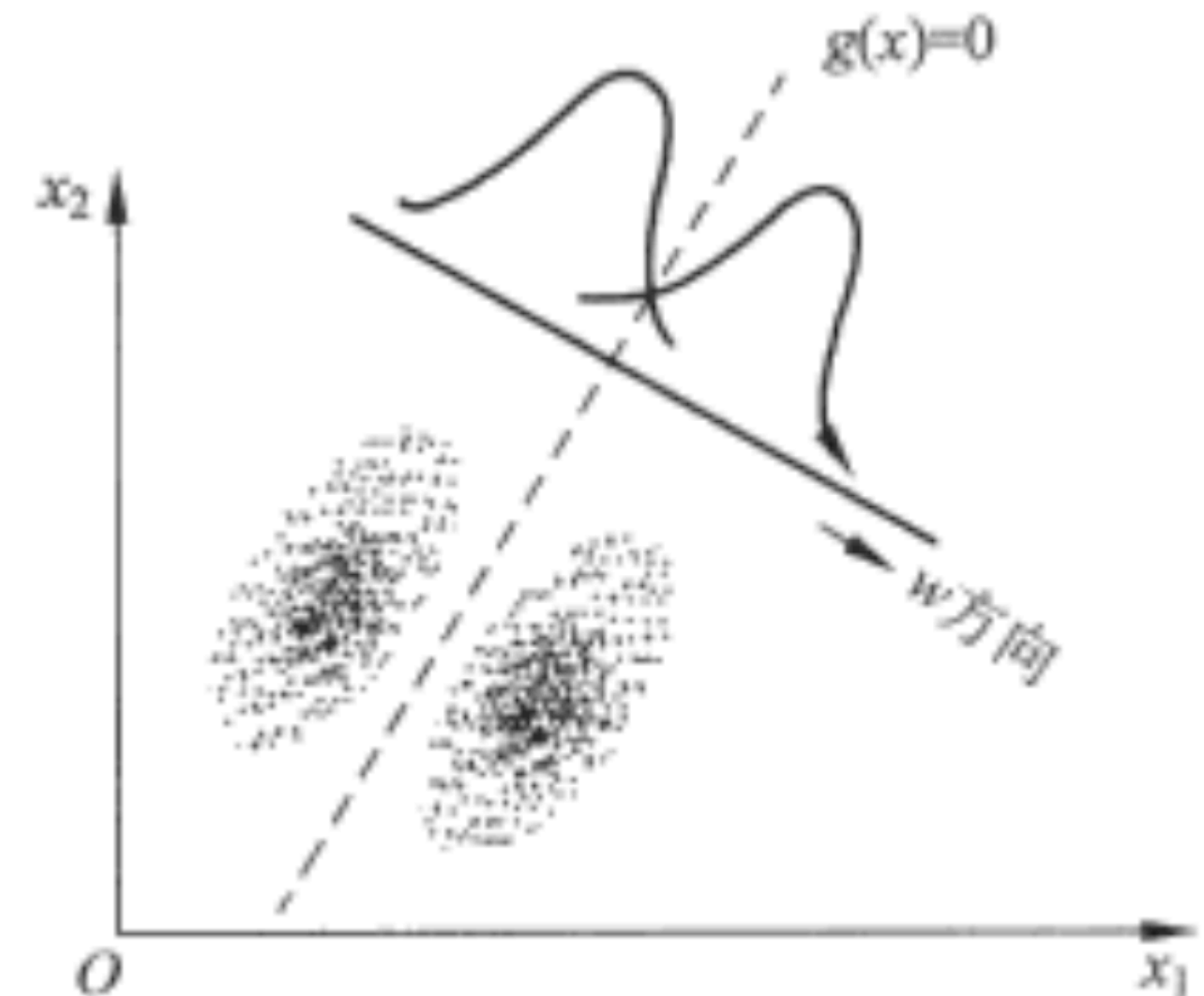
- 样本不服从正态分布时

- 可依据经验, 如 $w_0 = -\frac{1}{2}(\tilde{m}_1 + \tilde{m}_2)$ 或者 $w_0 = -\tilde{m}$ , ( $\tilde{m}$ 是所有样本在投影后的均值)

- 若 $g(\mathbf{x}) = \mathbf{w}^T \left( \mathbf{x} - \frac{1}{2} (\mathbf{m}_1 + \mathbf{m}_2) \right) \geq \log \frac{P(\omega_2)}{P(\omega_1)}$ , 则

$$\mathbf{x} \in \begin{cases} \omega_1 \\ \omega_2 \end{cases}$$

- 通过与两类均值投影的平分点相比较做分类决策



# 5.5 感知器

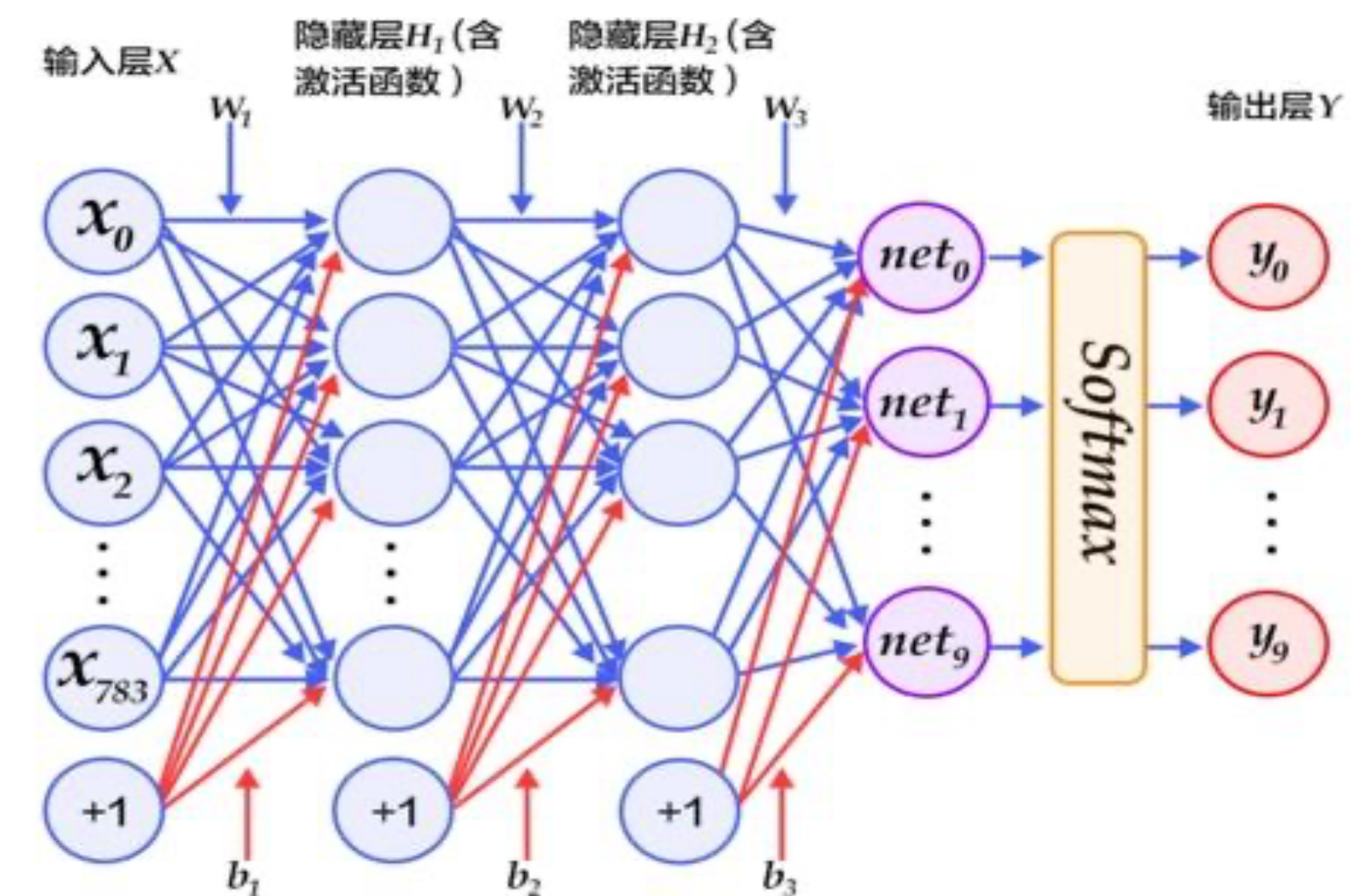
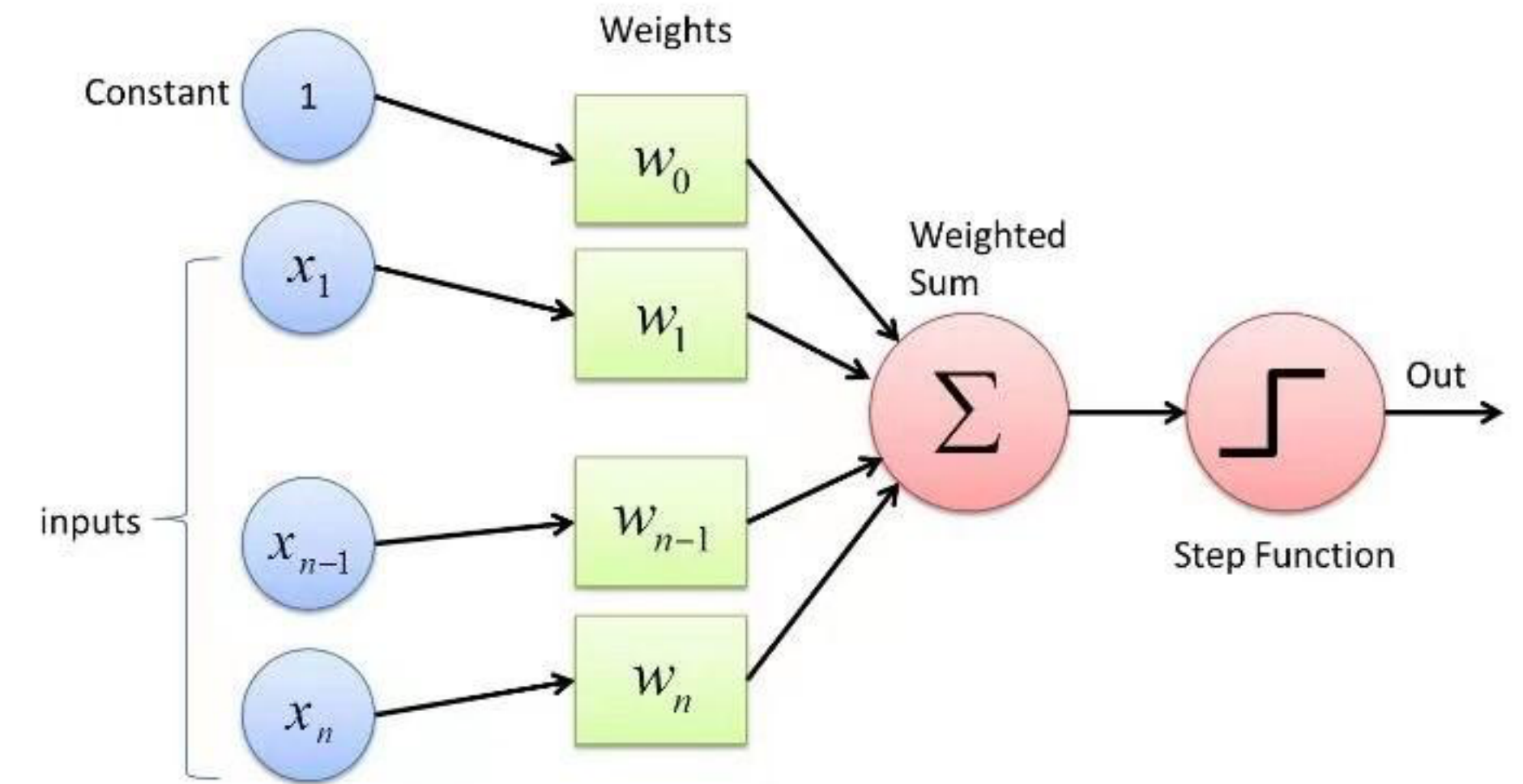
- **Fisher 线性判断：两步法**
  - 确定最优的投影方向
  - 在投影方向上确定分类阈值

- **感知器**

- 直接得到完整的线性判别函数

$$g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$$

- 多层感知器神经网络和各种深度学习方法的基础



# 5.5 感知器

- 线性判别函数的齐次简化

$$g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 \quad \rightarrow \quad g(\mathbf{y}) = \boldsymbol{\alpha}^T \mathbf{y}$$

$$\text{其中 } \mathbf{y} = [1, x_1, x_2, \dots, x_d]^T, \quad \boldsymbol{\alpha} = [w_0, w_1, w_2, \dots, w_d]^T$$

–  $\mathbf{y}$ 为增广的样本向量,  $\boldsymbol{\alpha}$ 为增广的权向量

- 决策规则 (两类 $\omega_1$ 、 $\omega_2$ )

$$\text{– } g(\mathbf{y}) \leq 0, \text{ 则 } \mathbf{y} \in \begin{cases} \omega_1 \\ \omega_2 \end{cases}$$

# 5.5 感知器

- 线性可分性

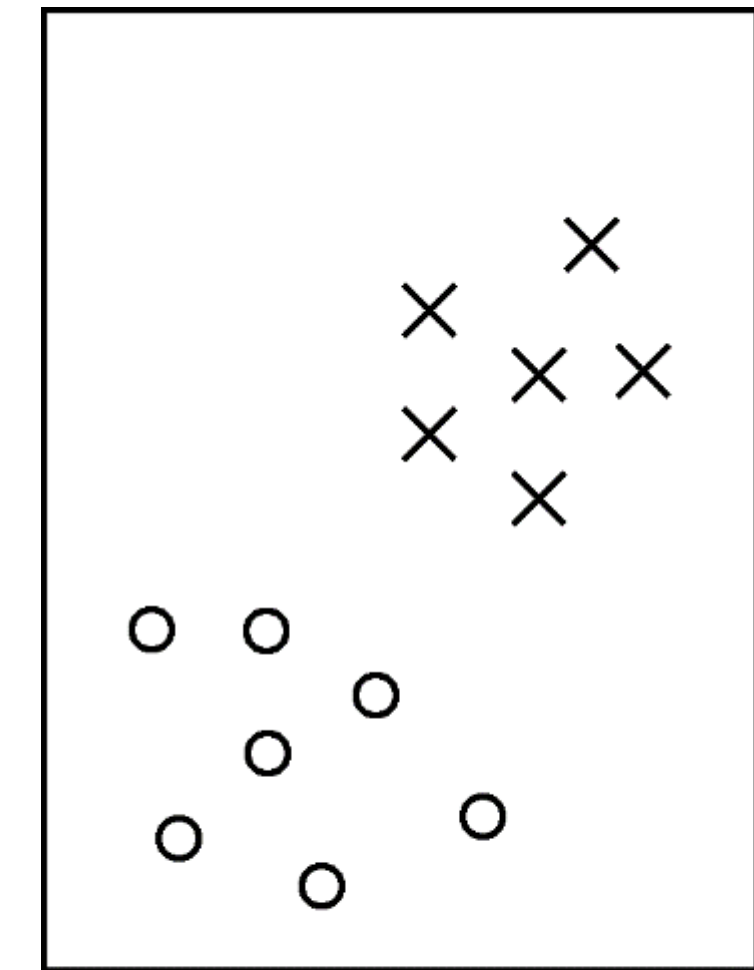
- 存在一个权向量 $\alpha$ ，使所有样本被正确分类，即

$$\exists \alpha, \forall i = 1, \dots, N, \begin{cases} \text{若 } y_i \in \omega_1, & \text{则 } \alpha^T y_i > 0 \\ \text{若 } y_i \in \omega_2, & \text{则 } \alpha^T y_i < 0 \end{cases}$$

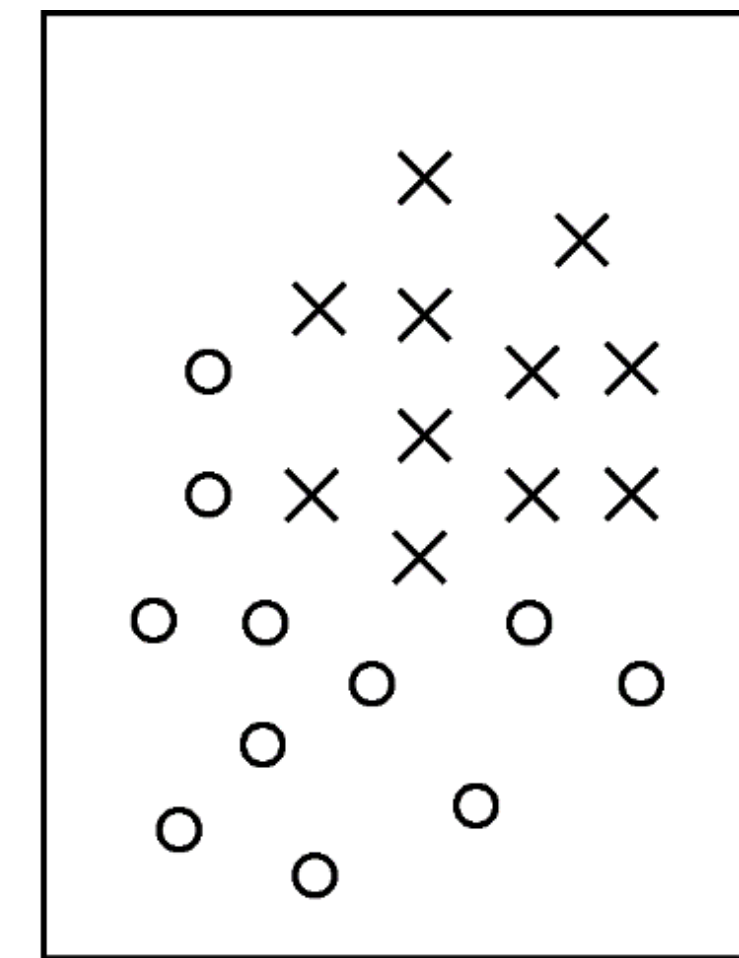
- 令 $y'_i = \begin{cases} y_i, & \text{若 } y_i \in \omega_1 \\ -y_i, & \text{若 } y_i \in \omega_2 \end{cases}$ ，则线性可分性变为：

$$\exists \alpha, \forall i = 1, \dots, N, \alpha^T y'_i > 0$$

- $y'$ : 规范化增广样本矩阵，仍然记作 $y$



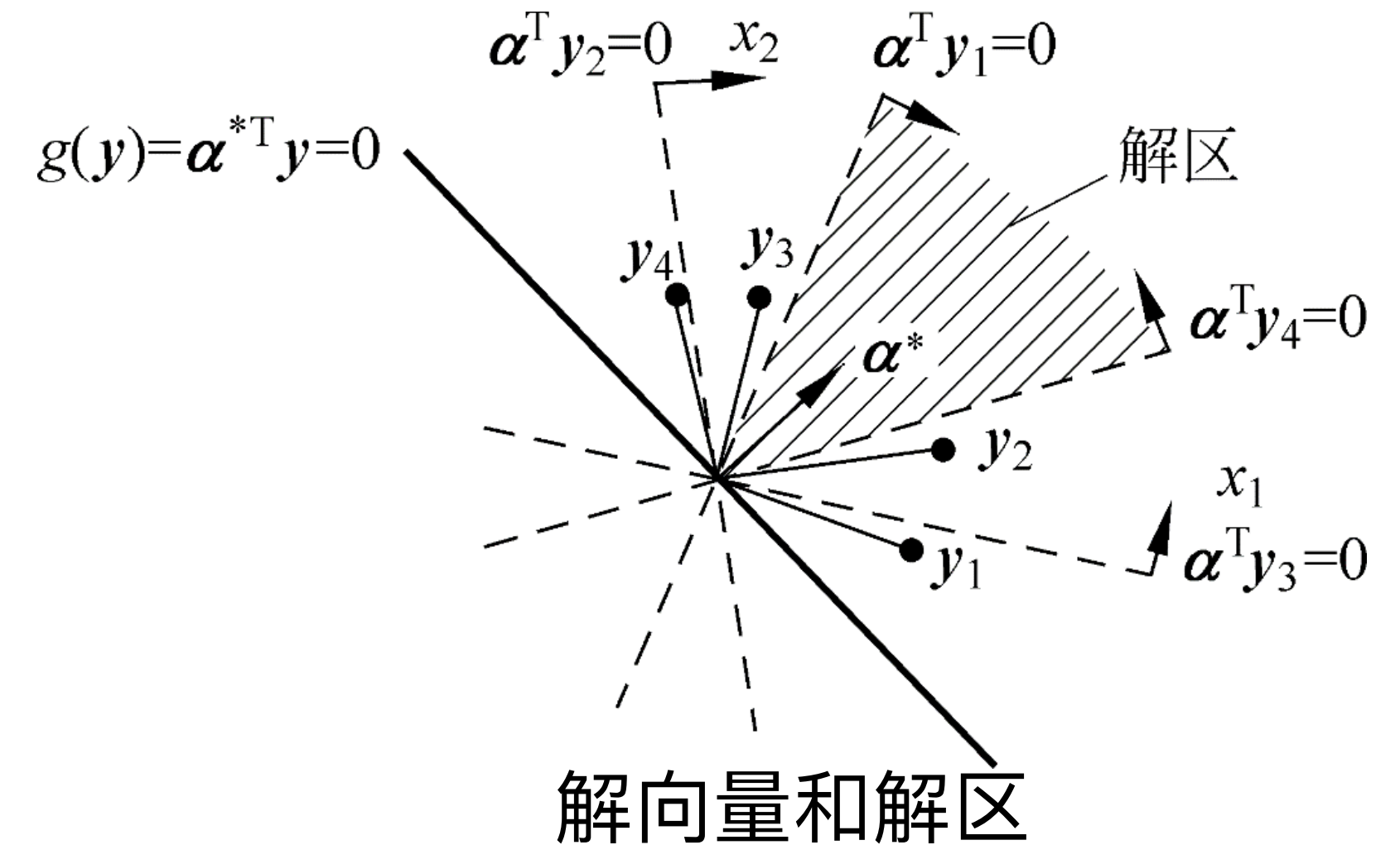
线性可分



线性不可分

# 5.5 感知器

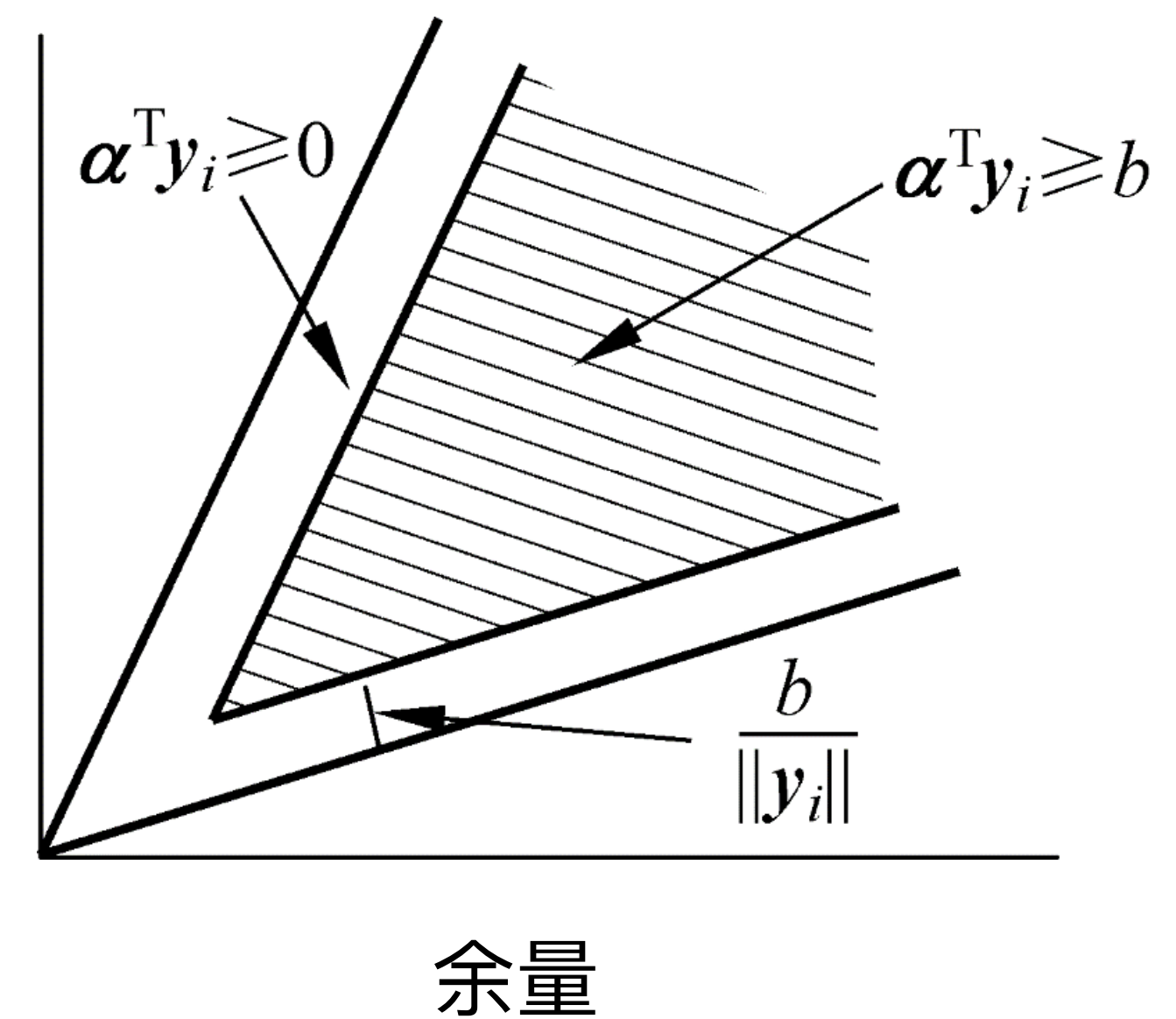
- 求解（线性可分情况）
  - 解向量 $\alpha^*$ ：满足 $\alpha^T y_i > 0, i = 1, 2, \dots, N$
  - 解区：权值空间中所有解向量组成的区域。
    - 对于样本 $y_i$ ,  $\alpha^T y_i = 0$ 定义了权值空间中一个超平面 $\hat{H}_i$ 。
    - 对于所有样本，解区就是每个样本对应超平面 $\hat{H}_i$ 的正侧的交集。



余量

# 5.5 感知器

- 求解 (线性可分情况)
  - 余量: 把解区向中间缩小, 不靠近边缘的解。  
即, 引入余量  $b > 0$ , 使得解向量满足  $\alpha^T y_i > b$ ,  
 $i = 1, 2, \dots, N$
  - 靠近边缘的解容易受到噪声、数值计算误差的影响, 解区中间更可靠。
  - B: 线性判别函数的平移



# 5.5 感知器

- 感知器准则函数

$$J_p(\boldsymbol{\alpha}) = \sum_{\boldsymbol{\alpha}^T \mathbf{y}_k \leq 0} (-\boldsymbol{\alpha}^T \mathbf{y}_k)$$

- 对所有错分样本的求和表示对错分样本的惩罚
- 当且仅当  $J_p(\boldsymbol{\alpha}^*) = \min J_p(\boldsymbol{\alpha}) = 0$  时,  $\boldsymbol{\alpha}^*$  是解向量

# 5.5 感知器

- 梯度下降法求解

$$J_p(\boldsymbol{\alpha}) = \sum_{\boldsymbol{\alpha}^T \mathbf{y}_k \leq 0} (-\boldsymbol{\alpha}^T \mathbf{y}_k)$$

- $\boldsymbol{\alpha}(t+1) = \boldsymbol{\alpha}(t) - \rho_t \nabla J_p(\boldsymbol{\alpha})$ ,  $\rho_t$ 为调整步长

- $\nabla J_p(\boldsymbol{\alpha}) = \frac{\partial J_p(\boldsymbol{\alpha})}{\partial \boldsymbol{\alpha}} = \sum_{\boldsymbol{\alpha}^T \mathbf{y}_k \leq 0} (-\mathbf{y}_k)$

- $\therefore \boldsymbol{\alpha}(t+1) = \boldsymbol{\alpha}(t) + \rho_t \sum_{\boldsymbol{\alpha}^T \mathbf{y}_k \leq 0} \mathbf{y}_k$ , 即每步迭代时把错分的样本按照某个系数加到权向量上。



# 5.5 感知器

- 梯度下降迭代算法步骤

- (1) 任意选择初始的权向量 $\alpha(0)$ , 置 $t = 0$

- (2) 每次考察一个样本 $y_j$ , 若 $\alpha(t)^T y_j \leq 0$ , 则 $\alpha(t+1) = \alpha(t) + y_j$ , 否则继续

- (3) 考察另一个样本, 重复(2), 直至对所有样本都有 $\alpha(t)^T y_j > 0$ , 即 $J_p(\alpha) = 0$

# 5.5 感知器

- 说明

- 若考虑余量 $b$ ，将错分判断条件变成 $\alpha(t)^T \mathbf{y}_j \leq b$ 即可
- 收敛性：对于线性可分的样本集，梯度下降的迭代算法经过有限次修正后一定会收敛到一个解向量 $\alpha^*$

- 可使用可变步长减少迭代步数：
$$\rho_i = \frac{|\alpha(k)^T \mathbf{y}_j|}{\|\mathbf{y}_j\|^2}$$

# 5.6 最小平方误差判别

- 考虑线性不可分情况

- $\alpha^T y_i > 0, i = 1, 2, \dots, N$ , 不等式组无法同时满足: 最小化错分样本数

- 形式化描述

- 把  $\alpha^T y_i > 0$  变成  $\alpha^T y_i = b_i > 0, i = 1, 2, \dots, N$

- 或矩阵形式  $Y\alpha = \mathbf{b}$ , 其中  $Y = \begin{bmatrix} y_1^T \\ \vdots \\ y_N^T \end{bmatrix} = \begin{bmatrix} y_{11} & \dots & y_1 \hat{d} \\ \vdots & \ddots & \vdots \\ y_{N1} & \dots & y_N \hat{d} \end{bmatrix}$ ,

- $\mathbf{b} = [b_1, b_2, \dots, b_N]^T$ ,  $\hat{d}$  是增广的样本向量的维数,  $\hat{d} = d + 1$

# 5.6 最小平方误差判别

- 求解

- 通常 $N > \hat{d}$ , 为矛盾方程组, 无法求得精确解

- 误差为 $e = Y\alpha - b$ , 可求得方程组的最小平方误差解, 即 $\alpha^*$ :  $\min J_s(\alpha)$

- $J_s(\alpha)$ 为最小平方误差 (MSE) 准则函数

$$J_s(\alpha) = \| Y\alpha - b \|^2 = \sum_{i=1}^N (\alpha^T y_i - b_i)^2$$

# 5.6 最小平方误差判别

- 求解

- 伪逆法求解

✓ 令  $\nabla J_s(\boldsymbol{\alpha}) = 2\mathbf{Y}^T(\mathbf{Y}\boldsymbol{\alpha} - \mathbf{b}) = 0$ , 可得  $\boldsymbol{\alpha}^* = (\mathbf{Y}^T\mathbf{Y})^{-1} \mathbf{Y}^T \mathbf{b} = \mathbf{Y}^+ \mathbf{b}$ ,

其中  $\mathbf{Y}^+ = (\mathbf{Y}^T\mathbf{Y})^{-1} \mathbf{Y}^T$  是长方矩阵  $\mathbf{Y}$  的伪逆

# 5.6 最小平方误差判别

- 求解

- 梯度下降法迭代求解

- ✓ 任意选择初始的权向量 $\alpha(0)$ , 置 $t = 0$

- ✓ 按照梯度下降的方式迭代更新权向量 $\alpha(t + 1) = \alpha(t) - \rho_t Y^T (Y\alpha - \mathbf{b})$ , 直到 $\nabla J_s(\alpha) \leq \xi$  或  $\|\alpha(t + 1) - \alpha(t)\| \leq \xi$  ( $\xi$ : 误差灵敏度)

- ✓ 也可通过单样本修正法调整权向量

- $\alpha(t + 1) = \alpha(t) + \rho_t (\mathbf{b}_k - \alpha(t)^T \mathbf{y}_k) \mathbf{y}_k$ , 其中,  $\mathbf{y}_k$ 是使得 $\alpha(t)^T \mathbf{y}_k \neq \mathbf{b}_k$

- 的样本: 最小均方根算法 (LMS 算法)、Widrow-Hoff算法

# 5.6 最小平方误差判别

- 求解

- 梯度下降法迭代求解

$\sqrt{b}$ 的选择

◆ 若  $\mathbf{b}_i = \begin{cases} \frac{N}{N_1} & \mathbf{y}_i \in \omega_1 \\ \frac{N}{N_2} & \mathbf{y}_i \in \omega_2 \end{cases}$ , 同类样本的  $\mathbf{b}_i$  相同, 则MSE的解等价于Fisher线性判别, 且

$w_0^* = -\mathbf{m}^T \mathbf{w}^*$ , 其中,  $N_1$ 、 $N_2$ 分别是第一和第二类样本数,  $N$ 为样本总数,

$$\mathbf{m} = \frac{1}{N}(N_1 \mathbf{m}_1 + N_2 \mathbf{m}_2)$$

◆ 若  $\mathbf{b}_i = 1 (\forall i = 1, \dots, N)$ , 则当  $N \rightarrow \infty$  时, MSE的解以最小平方误差逼近贝叶斯判别函数

$$g_0(\mathbf{x}) = P(\omega_1 | \mathbf{x}) - P(\omega_2 | \mathbf{x}), \text{ 即 } \boldsymbol{\alpha}^* \text{ 使得 } \varepsilon^2 = \int [\boldsymbol{\alpha}^T \mathbf{y} - g_0(\mathbf{x})]^2 p(\mathbf{x}) d\mathbf{x} \text{ 极小。}$$

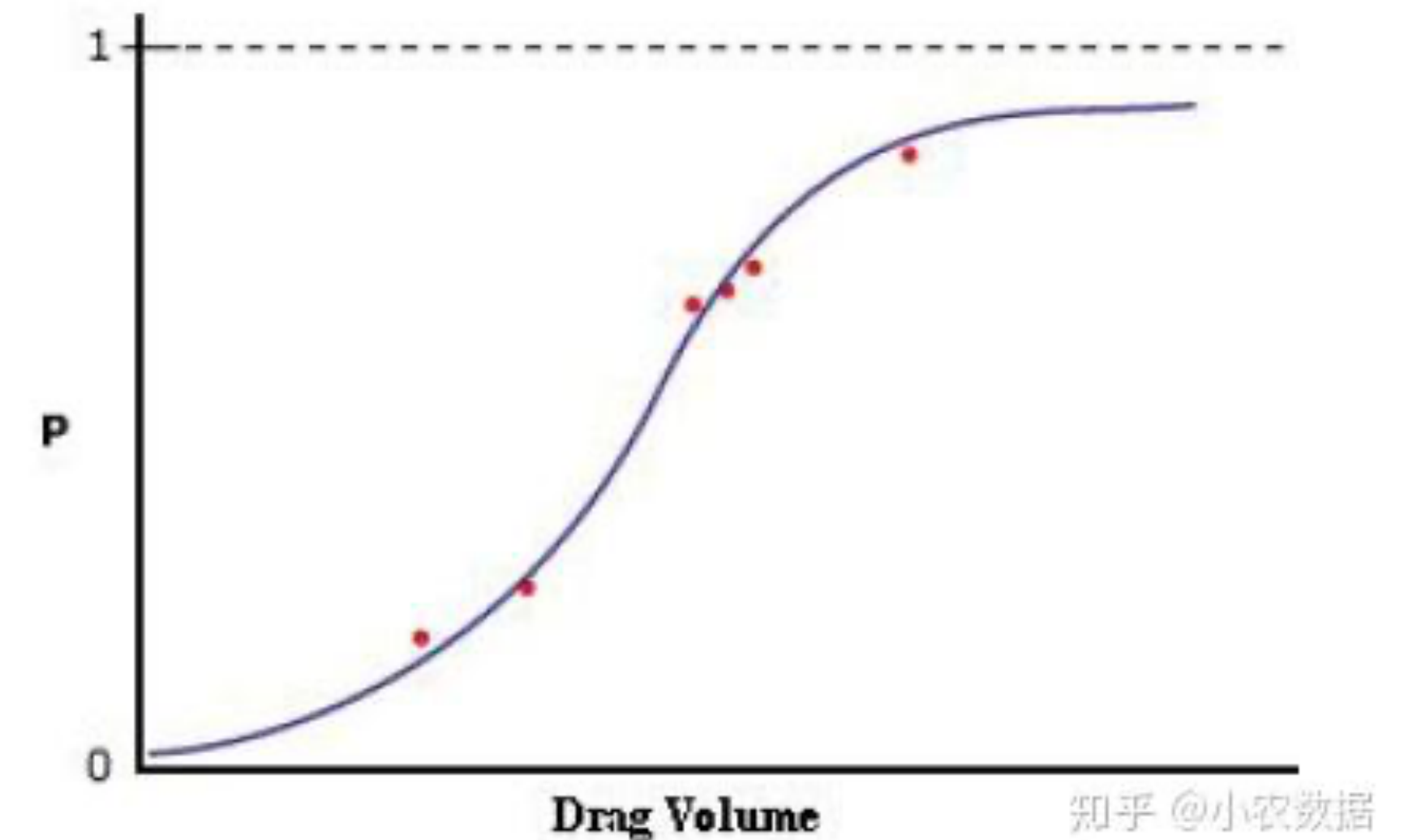
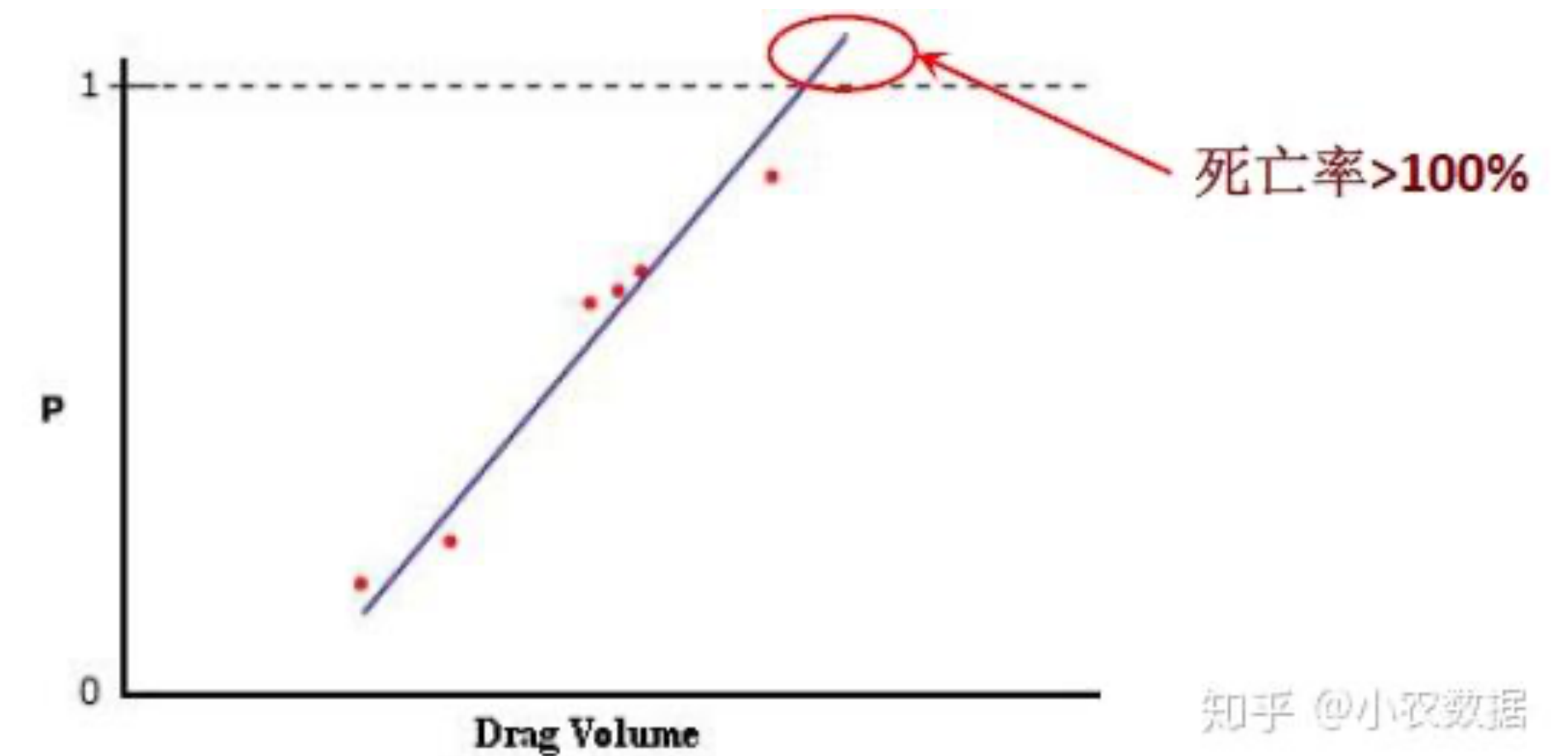
# 5.7 罗杰斯特回归 (logistic regression)

- 线性回归

- $y = \beta_0 + \beta_1 x_1 + \dots + \beta_m x_m + \varepsilon$

- 线性回归是**发散型方程式**，在许多应用上会有**不适合的情形**发生。

- 将线性回归的值压制在0~1之间，才能产生出合理的值。





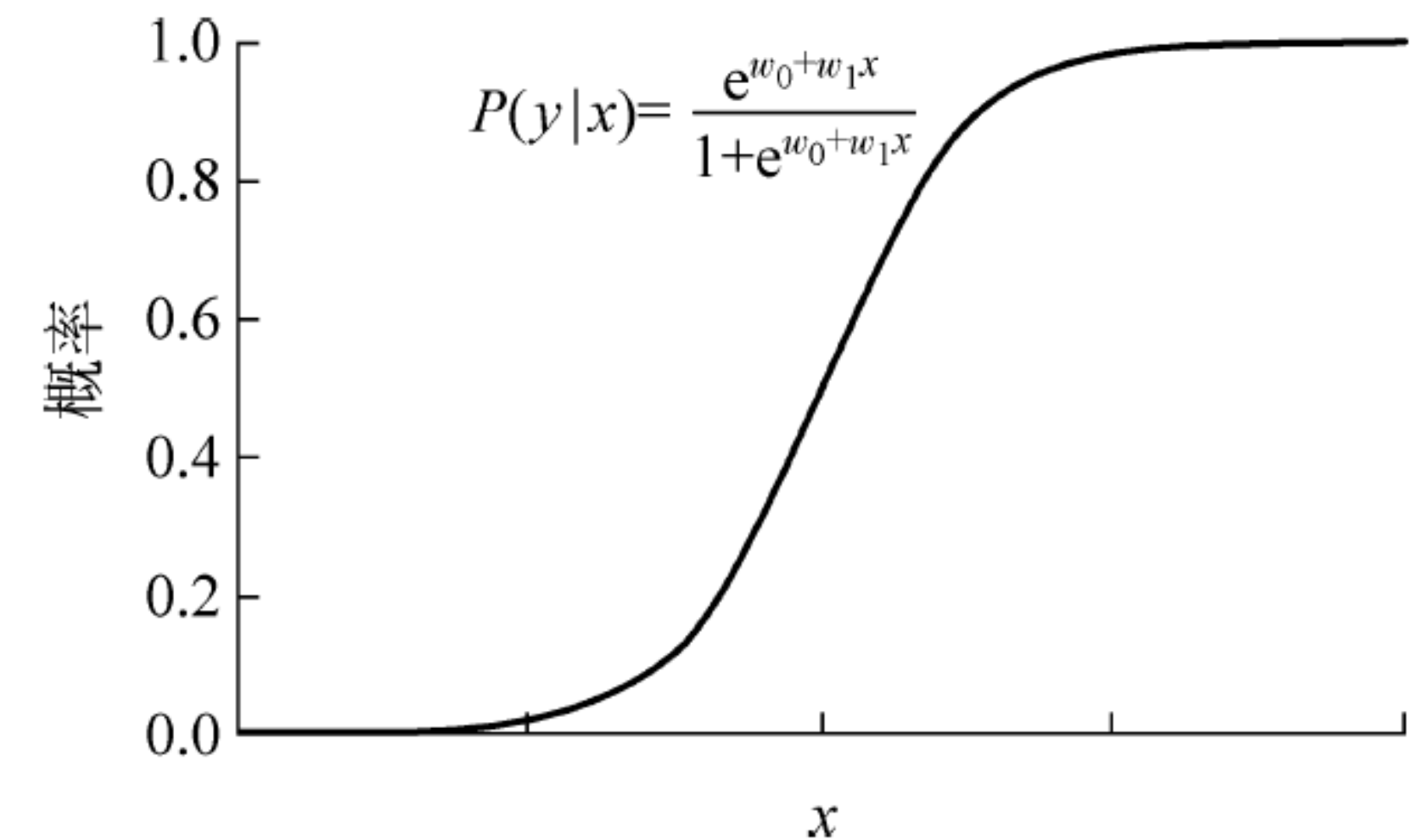
# 5.7 罗杰斯特回归 (logistic regression)

- 罗杰斯特函数 (logistic function)

- \_ 形式一:  $P(y = 1 | x) = P(y | x) = \frac{e^{w_0 + w_1 x}}{1 + e^{w_0 + w_1 x}}$

- \_ 形式二:  $\theta(s) = \frac{e^s}{1 + e^s} = \frac{1}{1 + e^{-s}}$ , 常用形式

- 在神经网络中称为Sigmoid函数



# 5.7 罗杰斯特回归 (logistic regression)

- 罗杰斯特回归 (logistic regression)

$$\text{logit}(\mathbf{x}) = \ln \left( \frac{P(y|\mathbf{x})}{1 - P(y|\mathbf{x})} \right) = \omega_0 + \omega_1 x_1 + \dots + \omega_m x_m$$

$$\triangleright P(y|\mathbf{x}) = \frac{e^{\omega_0 + \omega_1 x_1 + \dots + \omega_m x_m}}{1 + e^{\omega_0 + \omega_1 x_1 + \dots + \omega_m x_m}}$$

$$\triangleright \text{几率: } \frac{P(y|x)}{1 - P(y|x)} = e^{\omega_0 + \omega_1 x}, \text{ 医学中表示患病的可能性和不患病可能性之比}$$

$$\triangleright \text{对数几率: } \ln \left( \frac{P(y|x)}{1 - P(y|x)} \right) = \omega_0 + \omega_1 x$$

# 5.7 罗杰斯特回归 (logistic regression)

- 决策规则

- \_ 若  $\text{logit}(\mathbf{x}) \leq 0$ , 则  $\mathbf{x} \in \begin{cases} \omega_1 \\ \omega_2 \end{cases}$

- 决策规则学习算法: 最大似然法

- \_ 设共有  $N$  个独立的训练样本  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$ ,  $\mathbf{x}_j \in R^{d+1}$ ,

- $y_j \in \{-1, 1\}$ , 假设样本类别从某未知概率  $f(\mathbf{x})$  中产生, 以概率  $f(\mathbf{x})$  属于所关心

- 类别, 即  $P(y|\mathbf{x}) = \begin{cases} f(\mathbf{x}), & y = +1 \\ 1 - f(\mathbf{x}), & y = -1 \end{cases}$ , 用罗杰斯特函数  $h(\mathbf{x}) = \theta(\boldsymbol{\omega}^T \mathbf{x})$  估

- 计  $f(\mathbf{x})$ , 其中  $\boldsymbol{\omega}$  为罗杰斯特函数中待求参数组成的向量

# 5.7 罗杰斯特回归 (logistic regression)

- 决策规则学习算法：最大似然法

- 模型在样本上的似然函数：

$$P(y_j | \mathbf{x}_j) = \begin{cases} h(\mathbf{x}_j), & y_j = +1 \\ 1 - h(\mathbf{x}_j), & y_j = -1 \end{cases}$$

或者

$$l(h | (\mathbf{x}_j, y_j)) \triangleq P(y_j | \mathbf{x}_j, h) = \theta(y_j \mathbf{w}^T \mathbf{x}_j)$$

- 对于所有样本，模型的似然函数是

$$L(\mathbf{w}) = \prod_{j=1}^N P(y_j | \mathbf{x}_j) = \prod_{j=1}^N \theta(y_j \mathbf{w}^T \mathbf{x}_j)$$

# 5.7 罗杰斯特回归 (logistic regression)

## - 梯度下降法最优化目标函数

优化问题:

$$\begin{aligned}\min E(\mathbf{w}) &= -\frac{1}{N} \ln (L(\mathbf{w})) = -\frac{1}{N} \ln \left( \prod_{j=1}^N \theta \left( y_j \mathbf{w}^T \mathbf{x}_j \right) \right) \\ &= \frac{1}{N} \sum_{j=1}^N \ln \left( \frac{1}{\theta \left( y_j \mathbf{w}^T \mathbf{x}_j \right)} \right) = \frac{1}{N} \sum_{j=1}^N \ln \left( 1 + e^{-y_j \mathbf{w}^T \mathbf{x}_j} \right)\end{aligned}$$

# 5.7 罗杰斯特回归 (logistic regression)

## – 梯度下降法最优化目标函数

算法步骤:

✓ (1) 记时刻为  $k = 0$ , 初始化参数  $\omega(0)$

✓ (2) 计算目标函数的负梯度方向  $\nabla E = -\frac{1}{N} \sum_{j=1}^N \frac{y_j \mathbf{x}_j}{1 + e^{y_j \mathbf{w}(k)^T \mathbf{x}_j}}$

✓ 按步长 (学习率)  $\eta$  更新下一时刻参数  $\mathbf{w}(k+1) = \mathbf{w}(k) - \eta \nabla E$ , 检查是否达到终止条件, 如未达到, 令  $k = k + 1$ , 重新进行 (2)

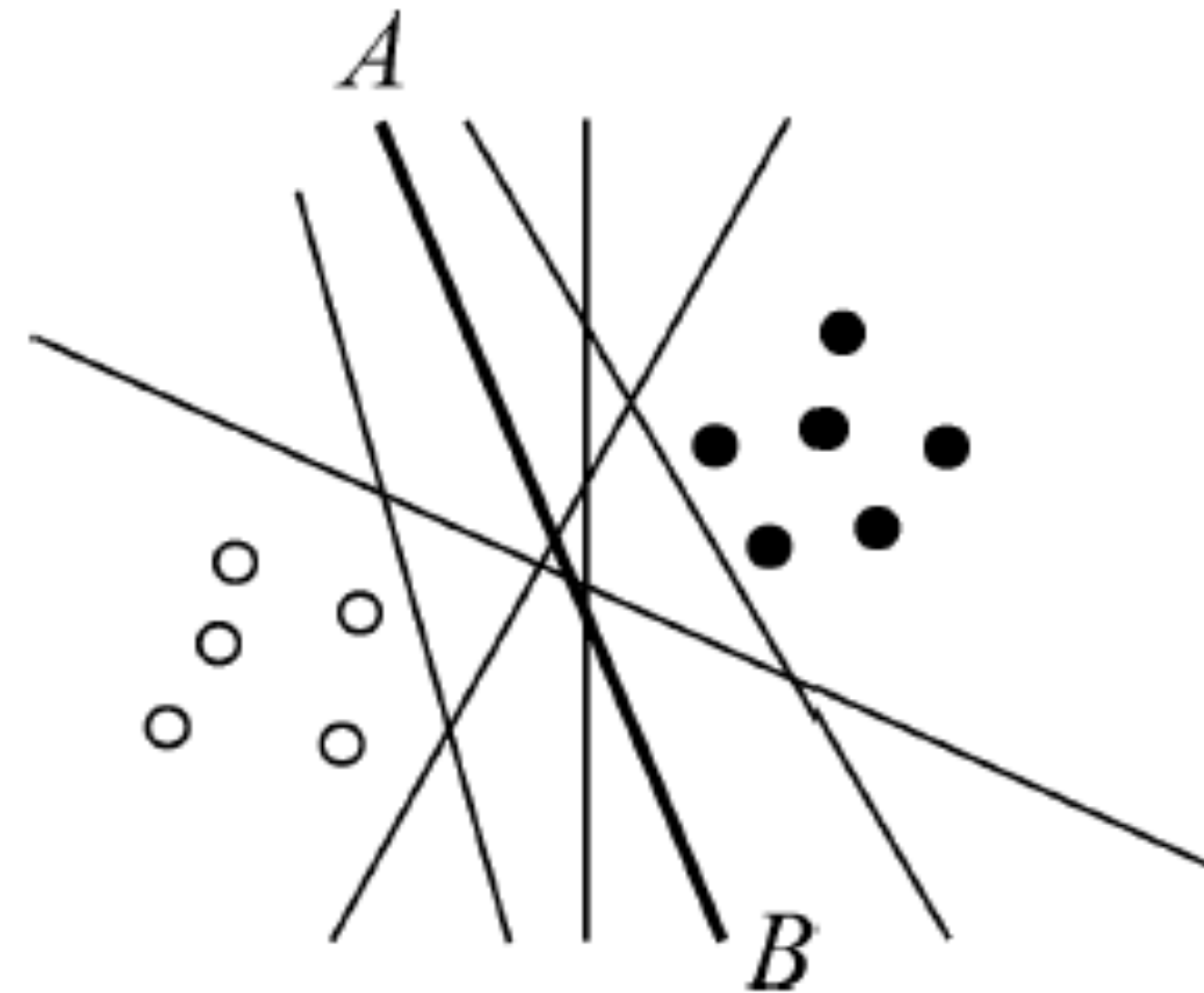
✓ (3) 算法停止, 输出得到的参数  $\omega$

终止条件可以是似然函数的梯度小于某个预设值, 训练过程不再有显著更新, 或迭代达到预设的上限等

# 5.8 最优分类超平面与线性支持向量机

- 问题提出

- 只要一个样本集线性可分，解区中的任何向量都是一个解向量
- 感知器算法采用不同的初值和迭代参数会得到不同的解
- 哪个解更好？



# 5.8.1 最优分类超平面

- 分类超平面

- 假设有训练样本集  $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)$ ,  $\mathbf{x}_i \in R^d$ ,  $y_i \in \{+1, -1\}$ , 其中每个样本是  $d$  维向量,  $y$  是类别标号,  $\omega_1$  类用  $+1$  表示,  $\omega_2$  类用  $-1$  表示。这些样本线性可分, 即存在超平面

$$g(\mathbf{x}) = (\mathbf{w} \cdot \mathbf{x}) + b = 0$$

把所有  $N$  个样本都没有错误地分开。

- 最优分类超平面

- 一个分类超平面, 如果它能将训练样本没有错误地分开, 且两类训练样本中离超平面最近的样本与超平面之间的距离是最大的, 则这个超平面称作**最优分类超平面** (optimal separating hyperplane), 简称**最优超平面** (optimal hyperplane)。

- ✓ 两类样本中离分类面最近的样本到分类面的距离称作**分类间隔** (margin)

- ✓ 最优超平面也称作**最大间隔超平面**



# 5.8.1 最优分类超平面

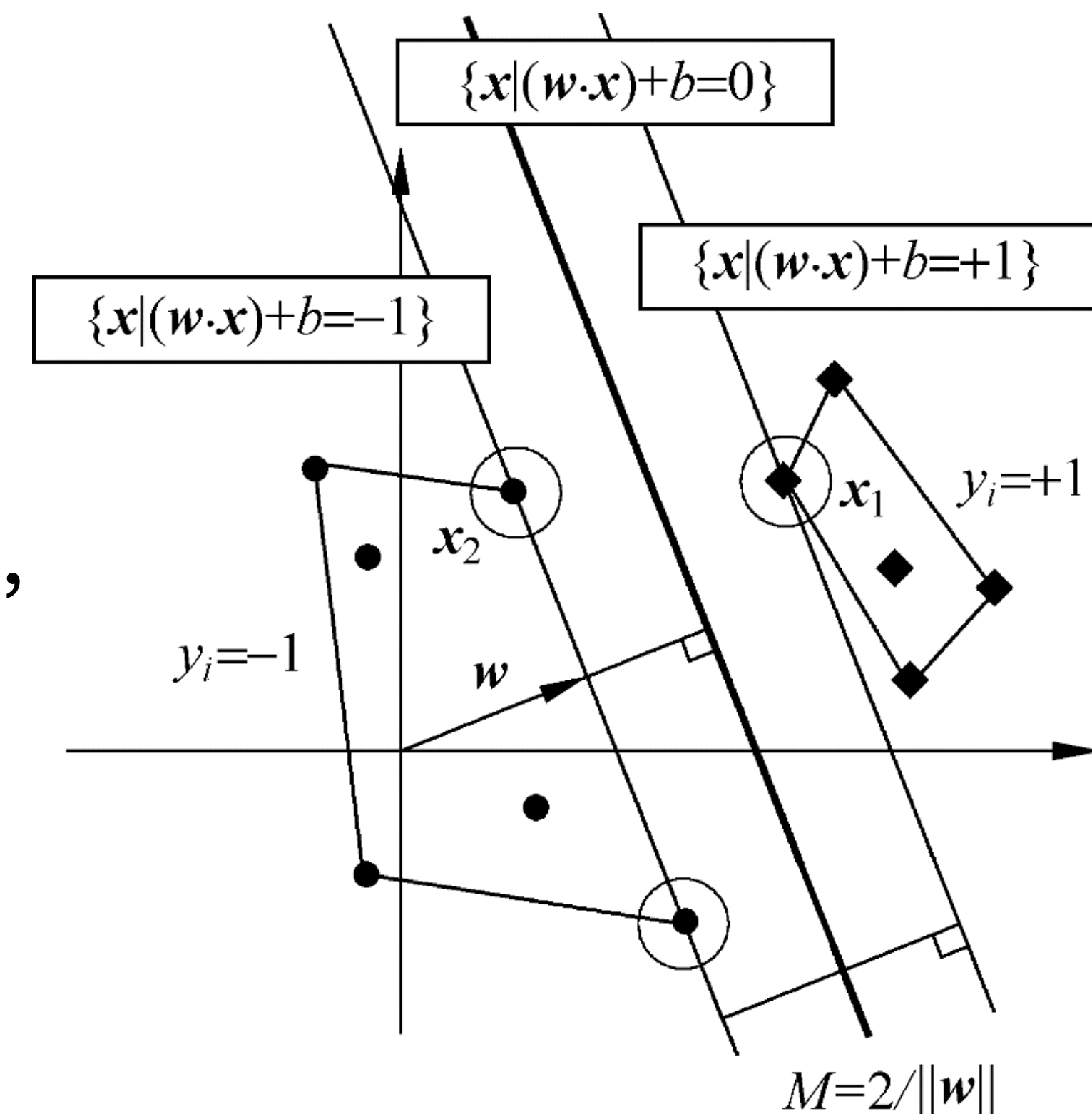
- 最优分类超平面

- 最优超平面定义的分类决策函数

$$f(\mathbf{x}) = \text{sgn}(g(\mathbf{x})) = \text{sgn}(\mathbf{w} \cdot \mathbf{x} + b)$$

- 规范化的分类超平面  $\begin{cases} (\mathbf{w} \cdot \mathbf{x}_i) + b \geq 1, y_i = +1 \\ (\mathbf{w} \cdot \mathbf{x}_i) + b \leq -1, y_i = -1 \end{cases}$

即  $y_i[(\mathbf{w} \cdot \mathbf{x}_i) + b] \geq 1, i = 1, 2, \dots, N$



# 5.8.1 最优分类超平面

- 最优分类超平面

- 求解最优超平面 ( $\mathbf{w}^*$ )

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 \quad \text{s.t.} \quad y_i \left[ (\mathbf{w} \cdot \mathbf{x}_i) + b \right] - 1 \geq 0, \quad i = 1, 2, \dots, N$$

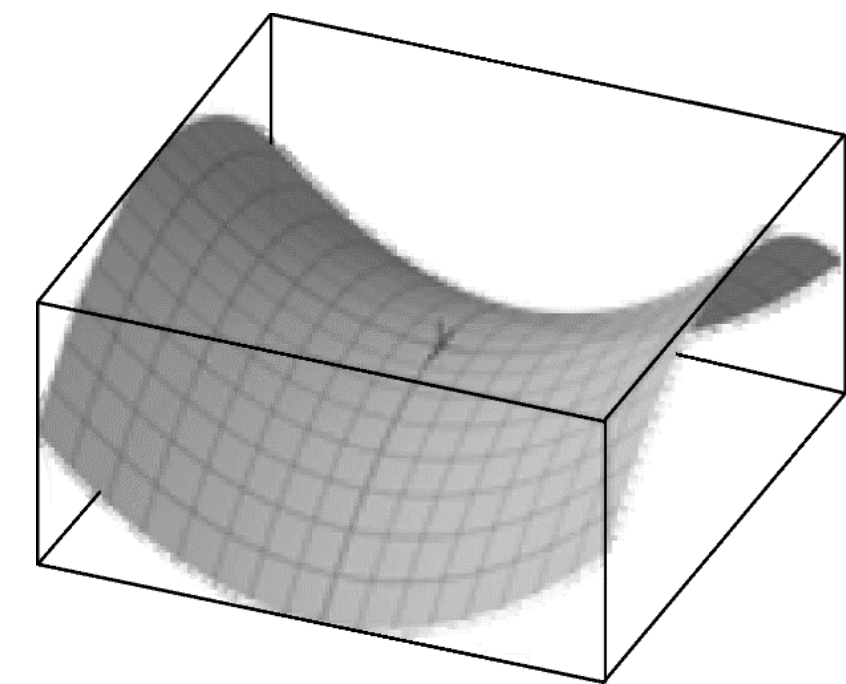
引入拉格朗日系数

$$\alpha_i \geq 0, \quad i = 1, \dots, N$$

等价问题

$$\min_{\mathbf{w}, b} \max_{\alpha} L(\mathbf{w}, b, \alpha) = \frac{1}{2} (\mathbf{w} \cdot \mathbf{w}) - \sum_{i=1}^N \alpha_i \{ y_i [ (\mathbf{w} \cdot \mathbf{x}_i) + b ] - 1 \}$$

最优解在  $L(\mathbf{w}, b, \alpha)$  的鞍点上取得，即  $L(\mathbf{w}, b, \alpha)$  对  $\mathbf{w}$  和  $b$  的偏导数均为 0



鞍点

# 5.8.1 最优分类超平面

- 最优分类超平面

- 求解最优超平面 ( $w^*$ )

最优解处:

$$w^* = \sum_{i=1}^N \alpha_i^* y_i x_i, \quad \text{且} \quad \sum_{i=1}^N y_i \alpha_i^* = 0$$

最优超平面的对偶问题 (the dual problem) :

$$\max_{\alpha} Q(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N \alpha_i \alpha_j y_i y_j (x_i \cdot x_j)$$

$$\text{s.t.} \quad \sum_{i=1}^N y_i \alpha_i = 0, \quad \alpha_i \geq 0, \quad i = 1, \dots, N$$

# 5.8.1 最优分类超平面

- 最优分类超平面

- 求解最优超平面 ( $\mathbf{w}^*$ )

通过对偶问题的解，可求出原问题的解：

$$\mathbf{w}^* = \sum_{i=1}^N \alpha_i^* y_i \mathbf{x}_i$$

$$f(\mathbf{x}) = \text{sgn}\{g(\mathbf{x})\} = \text{sgn}\{(\mathbf{w}^* \cdot \mathbf{x}) + b\} = \text{sgn}\left\{\sum_{i=1}^N \alpha_i^* y_i (\mathbf{x}_i \cdot \mathbf{x}) + b^*\right\}$$

# 5.8.1 最优分类超平面

- 最优分类超平面
  - 求解最优超平面 ( $b^*$ )
- 库恩-塔克 (Kuhn-Tucker) 条件：拉格朗日泛函的鞍点处满足

$$\alpha_i \{ y_i [(\mathbf{w} \cdot \mathbf{x}_i) + b] - 1 \} = 0, \quad i = 1, 2, \dots, N$$

等号成立的样本所对应的 $\alpha_i$ 才会大于0，求解 $b$ ：

$$y_i \left[ (\mathbf{w}^* \cdot \mathbf{x}_i) + b^* \right] - 1 = 0$$

# 5.8.2 线性不可分情况

- 样本集非线性可分

- 定义

对样本集  $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2) \cdots, (\mathbf{x}_N, y_N)$ ,  $\mathbf{x}_i \in R^d$ ,  $y_i \in \{+1, -1\}$ , 不等式

$$y_i[(\mathbf{w} \cdot \mathbf{x}_i) + b] - 1 \geq 0, \quad i = 1, 2, \dots, N$$

不可能对所有样本同时满足

- 策略

每个样本引入一个非负的松弛变量  $\xi_i$ , 约束条件变为

$$y_i[(\mathbf{w} \cdot \mathbf{x}_i) + b] - 1 + \xi_i \geq 0, \quad i = 1, 2, \dots, N$$

# 5.8.2 线性不可分情况

- 广义最优分类面

- 定义

$$\min_{\mathbf{w}, b, \xi_i} \frac{1}{2}(\mathbf{w} \cdot \mathbf{w}) + C \sum_{i=1}^N \xi_i$$

$$\text{s.t. } y_i \left[ (\mathbf{w} \cdot \mathbf{x}_i) + b \right] - 1 + \xi_i \geq 0, \quad i = 1, 2, \dots, N$$

$$\xi_i \geq 0, \quad i = 1, 2, \dots, N$$

# 5.8.2 线性不可分情况

- 广义最优分类面

- 求解

转化为拉格朗日泛函的鞍点问题

$$\begin{aligned} & \min_{\mathbf{w}, b, \xi_i} \max_{\alpha} L(\mathbf{w}, b, \alpha) \\ & = \frac{1}{2}(\mathbf{w} \cdot \mathbf{w}) + C \sum_{i=1}^N \xi_i - \sum_{i=1}^N \alpha_i \{ \gamma_i [(\mathbf{w} \cdot \mathbf{x}_i) + b] - 1 + \xi_i \} - \sum_{i=1}^N \gamma_i \xi_i \end{aligned}$$



# 5.8.2 线性不可分情况

- 广义最优分类面

- 求解

广义最优分类面的对偶优化问题

$$\max_{\alpha} Q(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j)$$

$$\text{s.t.} \quad \sum_{i=1}^N y_i \alpha_i = 0, \quad \text{且} \quad 0 \leq \alpha_i \leq C, \quad i = 1, 2, \dots, N$$

原问题的解向量满足：
$$\mathbf{w}^* = \sum_{i=1}^N \alpha_i^* y_i \mathbf{x}_i$$

广义最优分类面的判别函数：
$$f(\mathbf{x}) = \text{sgn}\{g(\mathbf{x})\} = \text{sgn}\{(\mathbf{w}^* \cdot \mathbf{x}) + b\} = \text{sgn}\left\{\sum_{i=1}^N \alpha_i^* y_i (\mathbf{x}_i \cdot \mathbf{x}) + b^*\right\}$$

# 5.8.2 线性不可分情况

- 广义最优分类面

- 求解

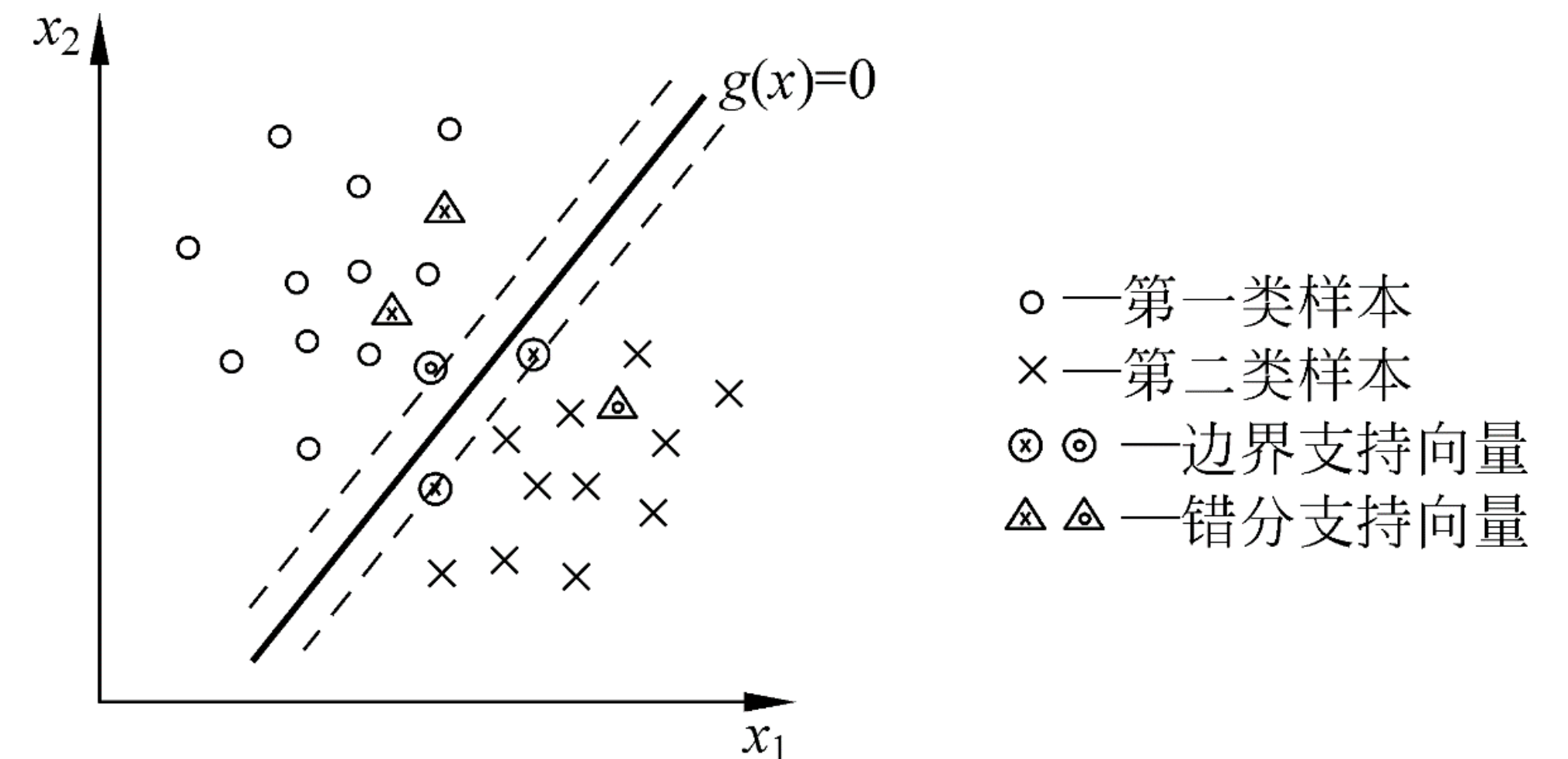
库恩-塔克条件：拉格朗日泛函的鞍点处满足

$$\alpha_i \{y_i[(\mathbf{w} \cdot \mathbf{x}_i) + b] - 1 + \xi_i\} = 0, \quad i = 1, 2, \dots, N$$

$$\gamma_i \xi_i = (C - \alpha_i) \xi_i = 0, \quad i = 1, 2, \dots, N$$

边界支持向量：  $0 < \alpha_i < C, \xi_i = 0$

错分支支持向量：  $\alpha_i = C, \xi_i > 0$



# 5.9 多类线性分类器

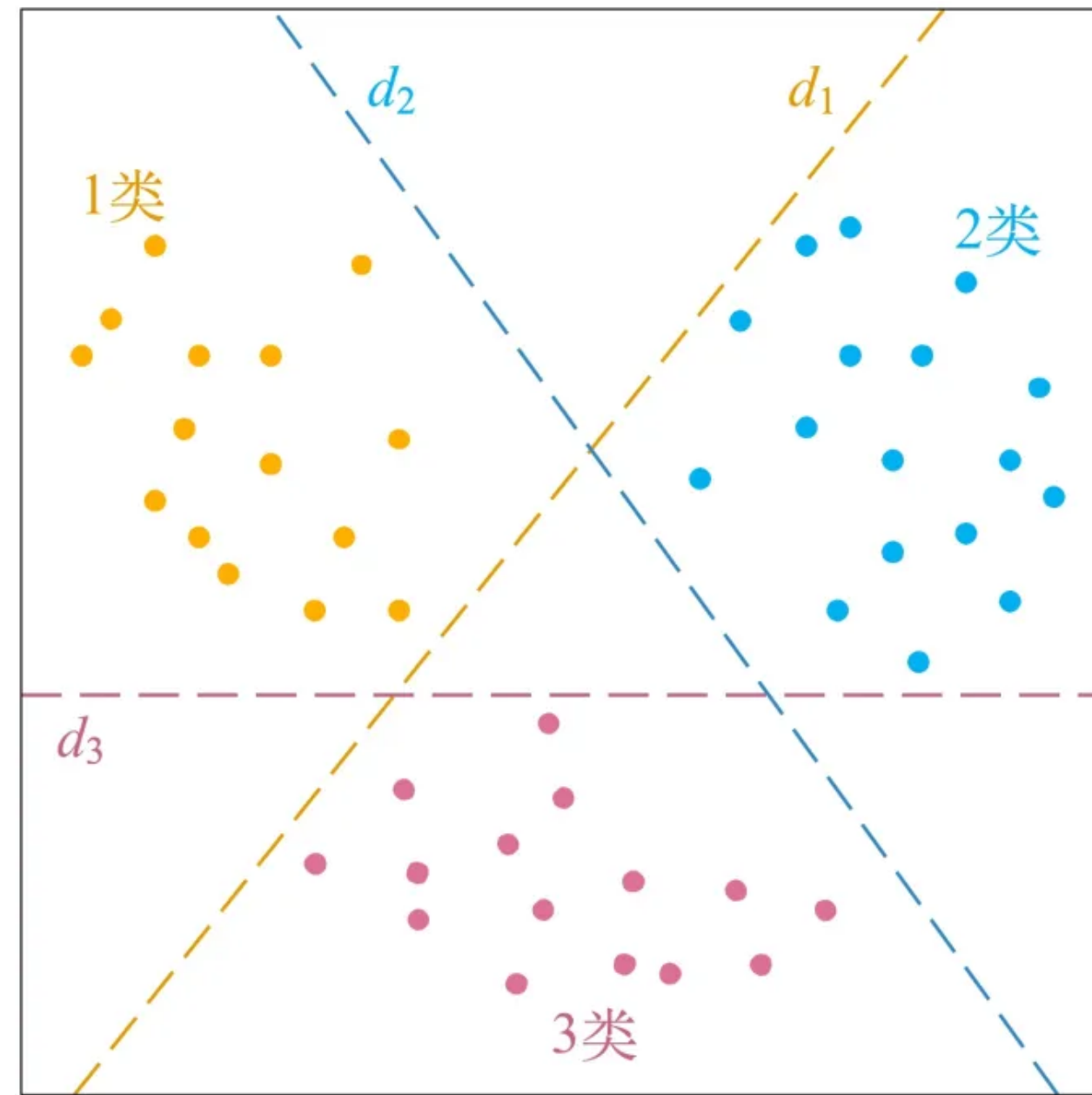
- 多分类问题

给定含 $N$ 个样本的训练集 $X = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$ , 其中 $K$ 维特征向量 $\mathbf{x}_n \in \mathbb{R}^K$ , 类标签 $y_n \in \{1, 2, \dots, M\}, n = 1, \dots, N$ 。训练集数据共 $M$ 个类。任务是找到决策函数 $y = f(\mathbf{x})$  (或者说一个规则)用于预测新数据的类别。

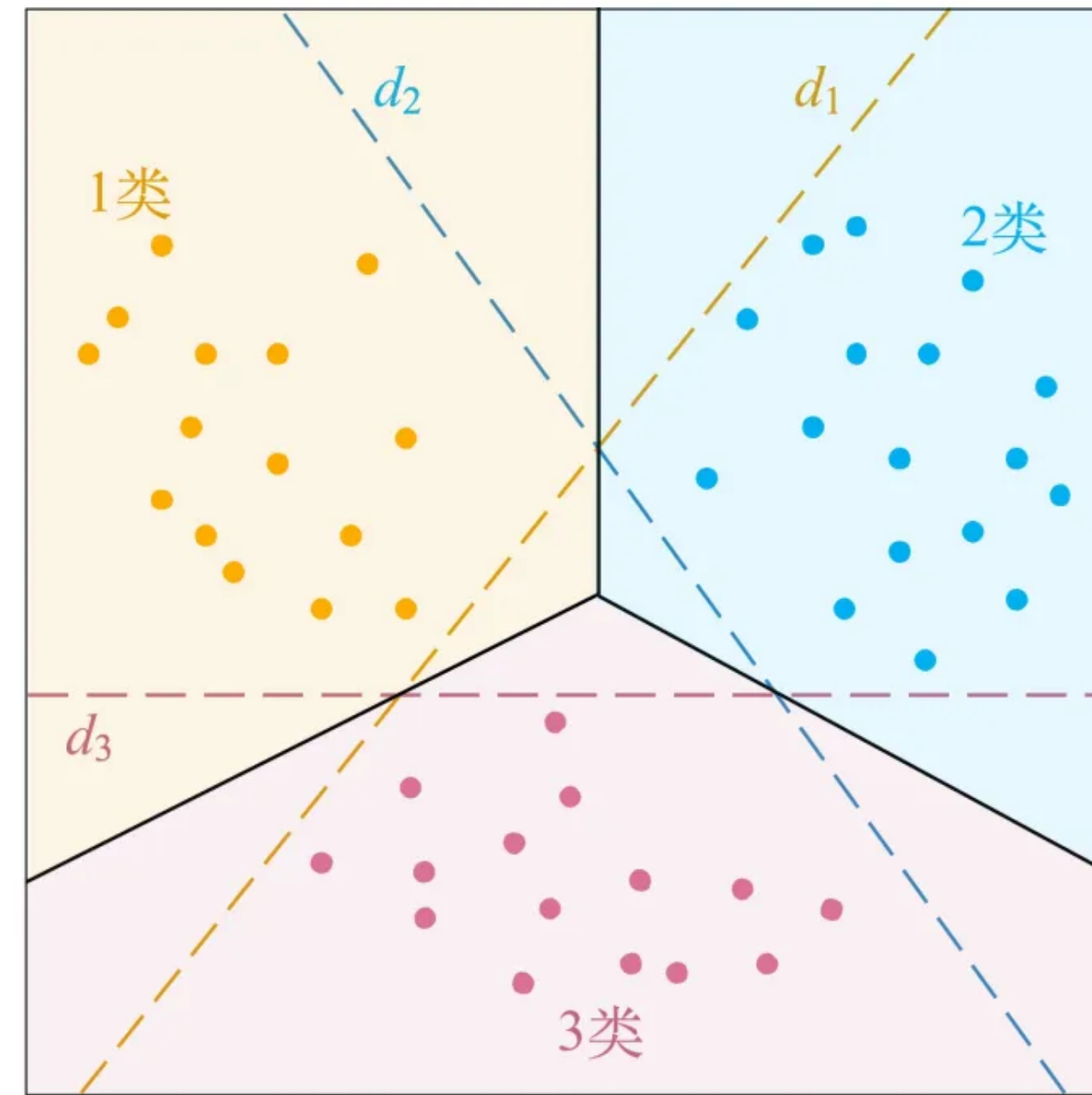
- 间接：将多分类问题分解为多个两类问题
- 直接：设计多类分类器

# 5.9.1 多个两类分类器的组合

- 方案一：一类对余类(one-against-all, one-against-the-rest)
  - $c$ 类转化为 $c - 1$ 个两类问题
  - 问题：训练样本不均衡；歧义区



(a)



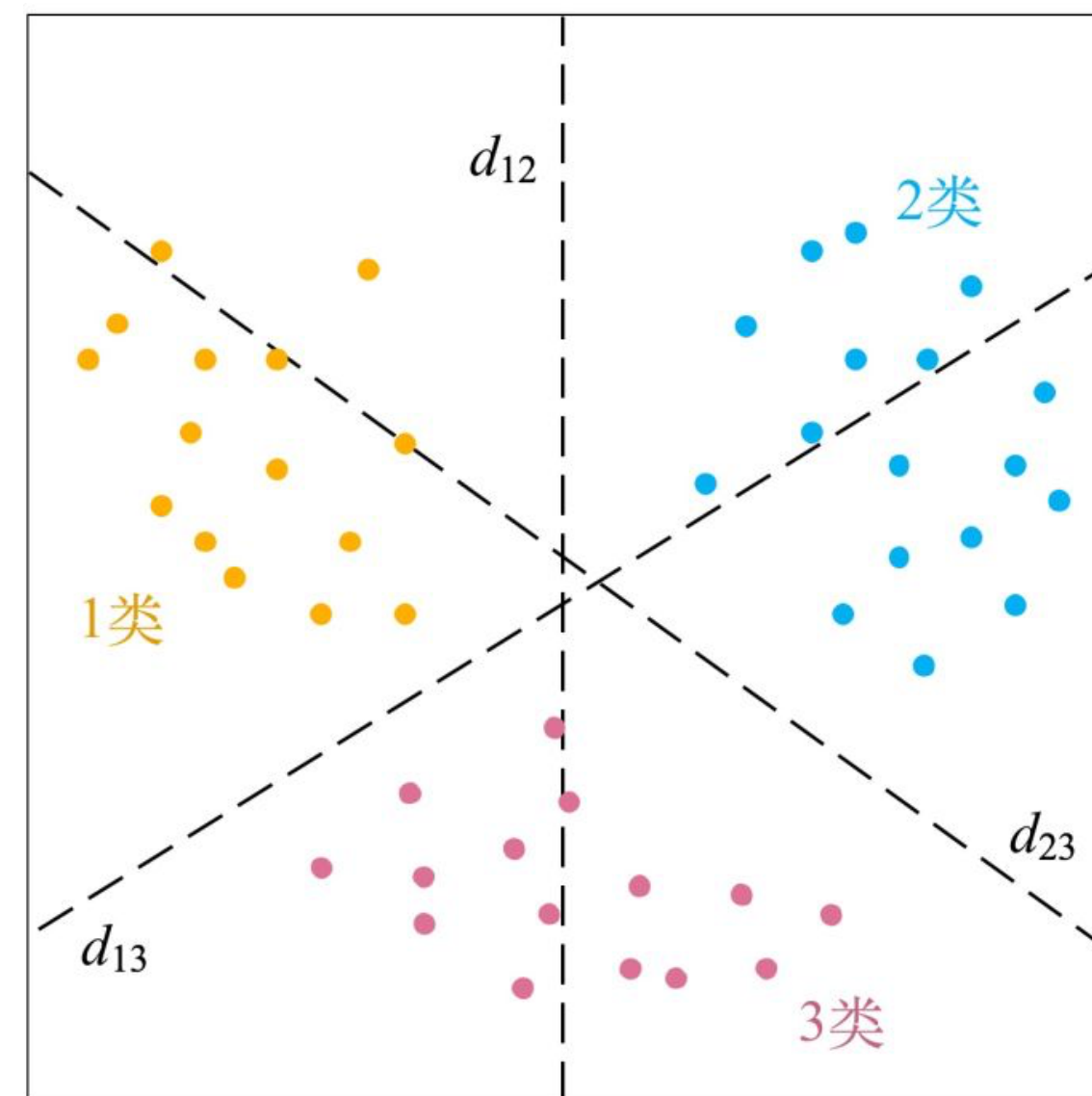
(b) 知乎 @RookieJ

# 5.9.1 多个两类分类器的组合

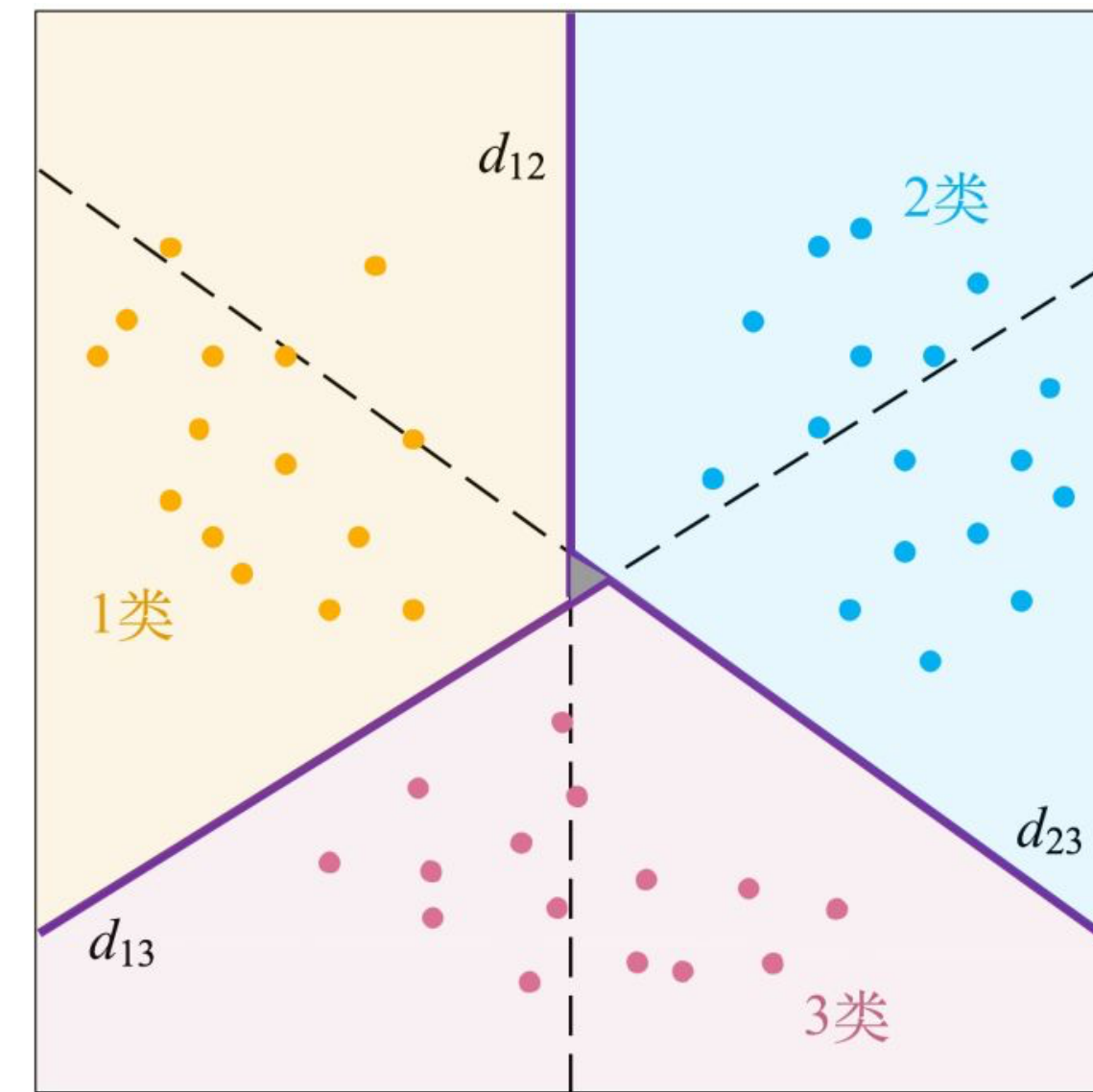
- 方案二：成对分类(one-against-one, pairwise classification)

– 问题：每两类构造一个分类器，则 $c$ 类需要 $\frac{c(c-1)}{2}$ 个两类分类器

– 问题：分类器多



(a)



(b)

知乎 @RookieJ

# 5.9.2 多类线性判别函数

- 定义

对 $c$ 类设计 $c$ 个判别函数

$$g_i(\mathbf{x}) = \mathbf{w}_i^T \mathbf{x} + w_{i0}, \quad i = 1, 2, \dots, c$$

取判别函数最大的类，即

$$\text{若 } g_i(\mathbf{x}) > g_j(\mathbf{x}), \quad \forall j \neq i, \text{ 则 } \mathbf{x} \in \omega_i$$

表示为增广向量形式： $g_i(\mathbf{x}) = \boldsymbol{\alpha}_i^T \mathbf{y}$ ,  $i = 1, 2, \dots, c$ , 其中,  $\boldsymbol{\alpha}_i = \begin{bmatrix} \mathbf{w}_i \\ w_{i0} \end{bmatrix}$  为增广权向量。

多类线性判别函数也称作多类线性机器，可记作 $L(\boldsymbol{\alpha}_1, \boldsymbol{\alpha}_2, \dots, \boldsymbol{\alpha}_c)$ 。

# 5.9.2 多类线性判别函数

- 求解：逐步修正法求解线性机器（多类线性可分）

(1) 任意选择初始的权向量 $\alpha_i(0)$ ,  $i = 1, 2, \dots, c$ , 置 $t = 0$

(2) 考察某个样本 $\mathbf{y}^k \in w_i$ , 若 $\alpha_i(t)^T \mathbf{y}^k > \alpha_j(t)^T \mathbf{y}^k$ , 则所有权向量不变; 若存在某个类 $j$ , 使 $\alpha_i(t)^T \mathbf{y}^k \leq \alpha_j(t)^T \mathbf{y}^k$ , 则选择 $\alpha_j(t)^T \mathbf{y}^k$ 最大的类别 $j$ , 对各类权值进行如下修正

$$\begin{cases} \alpha_i(t+1) = \alpha_i(t) + \rho_t \mathbf{y}^k \\ \alpha_j(t+1) = \alpha_j(t) - \rho_t \mathbf{y}^k \\ \alpha_l(t+1) = \alpha_l(t), \quad l \neq i, j \end{cases}$$

(3) 如果所有样本都分类正确, 则停止; 否则考查另一个样本, 重复 (2)