

# 分类模型

## 第一部分 线性分类模型

张朝晖

2018-2019学年 20180911-0912

## 1. 分类模块

### 产生式分类模型

#### A. 贝叶斯分类模型

### 判别式分类模型

#### 线性分类模型

- B. Fisher 判别分类
- C. 感知器分类模型
- D. 大间隔分类模型 (线性 *SVM*)

#### 非线性分类模型

- E. 核 *SVM* (非线性 *SVM*)
- F. 核 Fisher 判别分类
- G. 神经网络

### 其它分类模型

- H. *KNN* 分类模型
- I. 决策树分类模型
- J. *Logistic* 回归
- K. *Softmax* 回归

## 2. 聚类模块

- L. *K*-均值聚类
- M. 高斯混合聚类
- N. *DBSCAN* 聚类
- O. 层次聚类

## 3. 回归模块

- P. *KNN* 回归
- Q. 回归树
- R. 最小二乘线性回归
- S. 岭回归
- T. *LASSO* 回归

## 4. 集成学习

- U. *Bagging*
- V. 随机森林
- W. *Boosting*

## 5. 特征工程

- X. 主成分分析 (*PCA*)
- ...

## 6. 评价模块

混淆矩阵 (及其相关指标)、ROC 曲线、交叉验证

# 主要内容

A. 引言

B. Fisher线性判别分析

C. 感知器

D. 线性支持向量机(线性SVM)

## A. 引言

# 1. 有关约定

## 关键词:

- 线性运算、非线性运算
- 函数、方程
- 线性/非线性
- 特征空间及其划分
- 决策域、决策边界(决策面、分类边界)
- 产生式分类模型、判别式(鉴别式)分类模型
- 线性/非线性
  - 分类器、分类模型、分类边界、判别函数
- 两类别分类、多类别分类

## 2.基于**模型**的分类器设计 -- 代表模型：贝叶斯分类模型

### 基本思路：

(1)基于训练集，估计描述每个类别的**概率密度(参数)**

(2)观测样本的类别决策

训练样本  $\rightarrow$  **估计 $p(x | \omega_j)$**   $\rightarrow$  贝叶斯判决

### 如何估计了类条件概率密度函数

#### (1)参数估计法

需明确类条件概率密度函数形式；

参数估计需要大量样本

#### (2)非参数估计法

不足： $p(x | \omega_j)$ 的估计需要大量样本；

随特征空间维数增加，占用大量存储空间

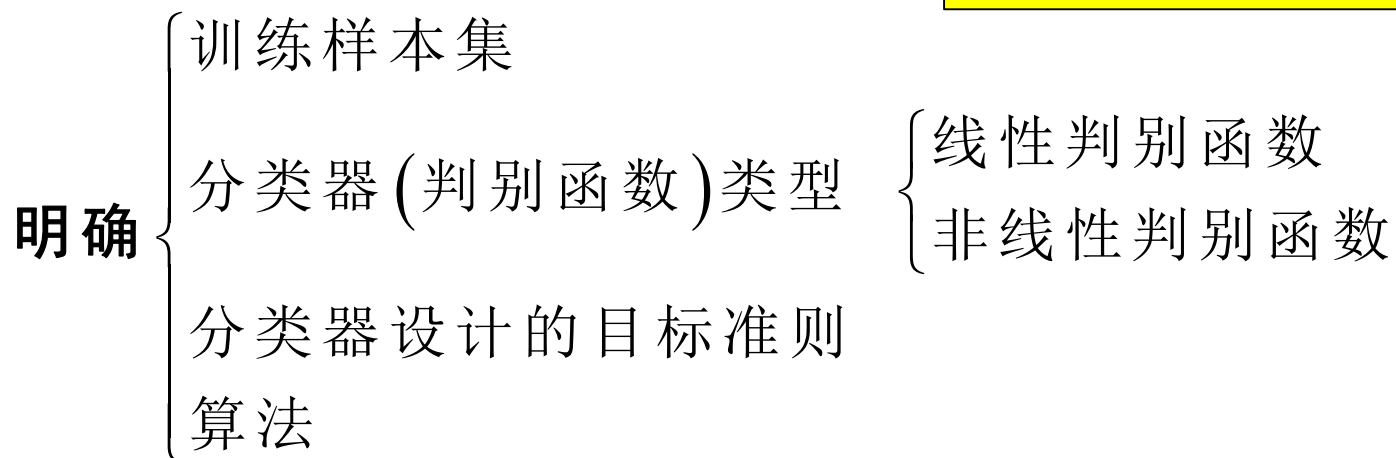
**$\Rightarrow$ 贝叶斯决策尽管最优，但实现困难**

### 3. 基于样本(数据)的直接分类器设计

例：鉴别式分类模型，由训练样本集，直接确定决策域划分

#### 基本思路

这里我们首先关注利用样本数据  
设计鉴别式分类模型!!!



——→ 利用样本，估计判别函数的未知参数  
使所选目标准则最优

分类器设计  $\Leftarrow$  **判别函数  
参数估计**  $\Rightarrow$  准则函数极值解求取

## 4. 鉴别式分类模型

线性分类模型

非线性分类模型

先从简单情况开始：  
面向两类别分类的“线性分类模型”

线性分类模型的分类边界是关于 $\mathbf{x}$ 的线性方程.

### 线性判别函数

$$\left\{ \begin{array}{l} \text{对于多类: } g_i(\mathbf{x}) = \mathbf{w}_i^T \mathbf{x} + \omega_{i0} \quad i = 1, \dots, c \\ \quad \quad \quad c \text{ 个判别函数} \\ \text{对于两类: } g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + \omega_0 \quad 1 \text{ 个判别函数} \end{array} \right.$$

### 两类问题 $c = 2$

$$\left\{ \begin{array}{l} d \text{ 维特征向量 (样本向量)} \quad \mathbf{x} = [x_1 \quad x_2 \quad \cdots \quad x_d]^T \\ \text{判别函数} \quad g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + \omega_0 \quad [g(\mathbf{x}) = g_1(\mathbf{x}) - g_2(\mathbf{x})] \\ \text{决策边界} \quad \mathbf{w}^T \mathbf{x} + \omega_0 = 0 \\ \quad \quad \quad \text{法向量 (权向量)} \quad \mathbf{w} = [w_1 \quad w_2 \quad \cdots \quad w_d]^T \\ \quad \quad \quad \text{阈值权} \quad \omega_0 \\ \text{决策域: 第1类决策域、第2类决策域} \end{array} \right.$$



**两类别分类模型**--决策边界 $H$  (分类超平面)的确定:

$$g(\mathbf{x})=0 \quad (g(\mathbf{x})=\mathbf{w}^T \mathbf{x}_1 + w_0)$$

对于 $\forall \mathbf{x}_1, \mathbf{x}_2 \in H$  并且  $\mathbf{x}_1 \neq \mathbf{x}_2$  有  $g(\mathbf{x}_1)=0, g(\mathbf{x}_2)=0$

$$\text{即 } \mathbf{w}^T \mathbf{x}_1 + w_0 = \mathbf{w}^T \mathbf{x}_2 + w_0$$

$$\longrightarrow \mathbf{w}^T (\mathbf{x}_1 - \mathbf{x}_2) = 0$$

**权向量 $\mathbf{w}$ 与超平面 $H$ 的任意向量正交** ( $\mathbf{w}$ 垂直于 $H$ )

**权向量 $\mathbf{w}$ 是超平面 $H$ 的法向量**

对于 $\mathbf{x} \in \mathcal{R}_1$ , 有 $g(\mathbf{x}) > 0$ , 所以 **决策面 $H$ 的法向量指向 $\mathcal{R}_1$**

$$\longrightarrow \begin{cases} \text{决策面 } H \text{ 的正侧: } \mathcal{R}_1 \\ \text{决策面 } H \text{ 的负侧: } \mathcal{R}_2 \end{cases}$$

# 两类别分类模型

## 决策规则

对于观测 $\mathbf{x}$   $\begin{cases} \text{若 } g(\mathbf{x}) > 0, & \text{则决策 } \mathbf{x} \text{ 为 } \omega_1 \text{ 类} \\ \text{若 } g(\mathbf{x}) < 0, & \text{则决策 } \mathbf{x} \text{ 为 } \omega_2 \text{ 类} \\ \text{若 } g(\mathbf{x}) = 0, & \text{则拒绝, 或结合实际问题分到某1类} \end{cases}$

**决策面 (决策边界)  $H$**        $g(\mathbf{x}) = 0$  (超平面)

**决策域**  $\begin{cases} \mathcal{R}_1 = \{\mathbf{x} : \mathbf{x} \in \mathcal{R} \text{ 且 } g(\mathbf{x}) > 0\} \\ \mathcal{R}_2 = \{\mathbf{x} : \mathbf{x} \in \mathcal{R} \text{ 且 } g(\mathbf{x}) < 0\} \end{cases}$

可将判别函数 $g(x)$ 视为观测 $x$ 到决策面 $H$ 距离的代数度量。

判别函数  $g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + \omega_0$ 

三角形三边关系  $\mathbf{x} = \mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|}$

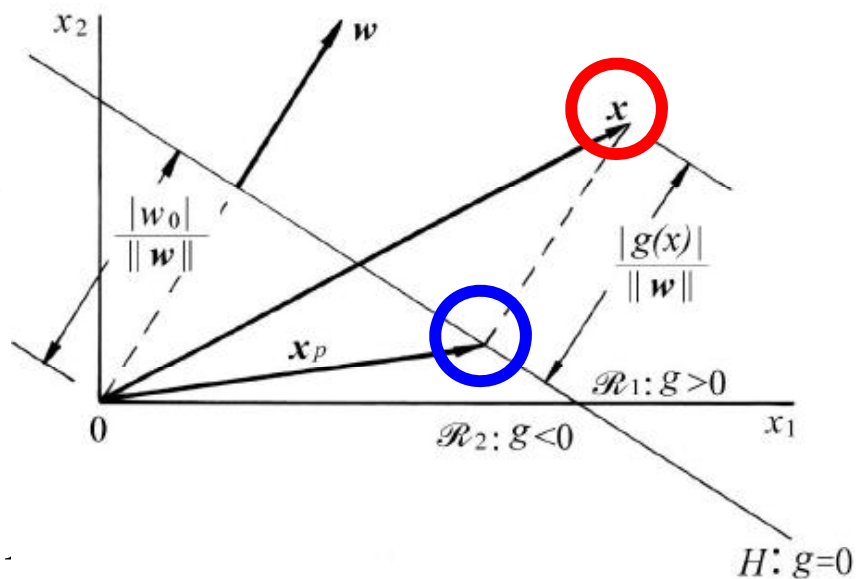
其中：

$\mathbf{x}_p$   $\mathbf{x}$ 在 $H$ 的投影向量,  $\mathbf{x}_p \in H$

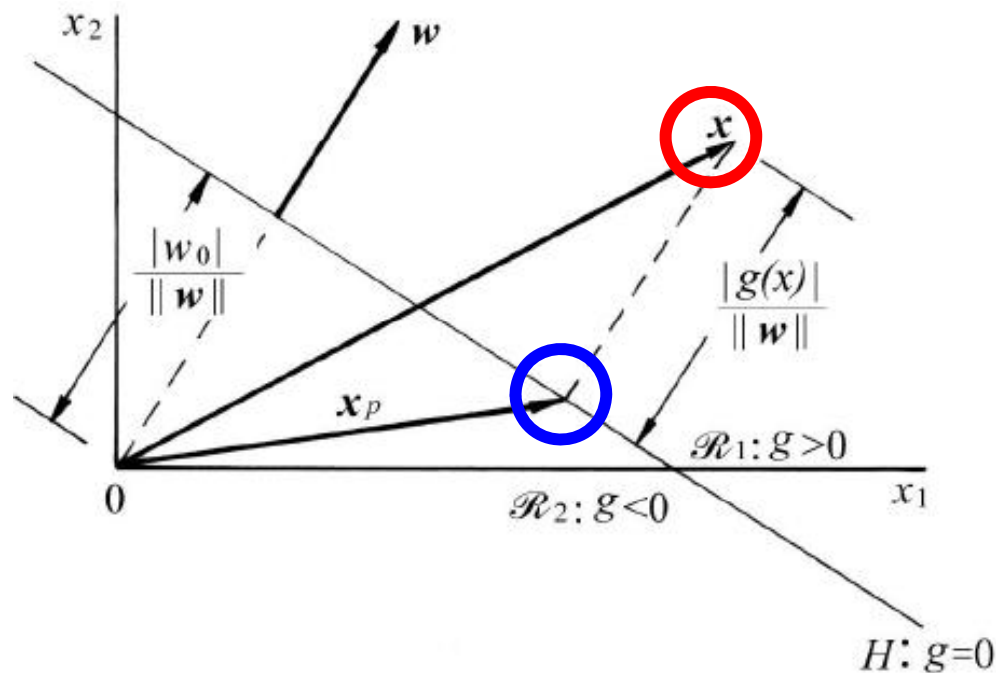
$$\mathbf{g}(\mathbf{x}_p) = 0 = \mathbf{w}^T \mathbf{x}_p + \omega_0$$

## $x$ 点到 $H$ 面的代数垂直距离

若 $r > 0$ ,则 $\mathbf{x}$ 在 $\mathbf{H}$ 正侧。

$$\frac{\boldsymbol{w}}{\|\boldsymbol{w}\|} \quad \boldsymbol{w} \text{方向的单位向量}$$


$$\begin{cases} g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + \omega_0 \\ g(\mathbf{x}_p) = 0 = \mathbf{w}^T \mathbf{x}_p + \omega_0 \\ \mathbf{x} = \mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|} \end{cases}$$



$$\Rightarrow g(\mathbf{x}) = \mathbf{w}^T \left( \mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|} \right) + \omega_0$$

$$= \mathbf{w}^T \mathbf{x}_p + r \mathbf{w}^T \frac{\mathbf{w}}{\|\mathbf{w}\|} + \omega_0 = g(\mathbf{x}_p) + r \|\mathbf{w}\| = r \|\mathbf{w}\|$$

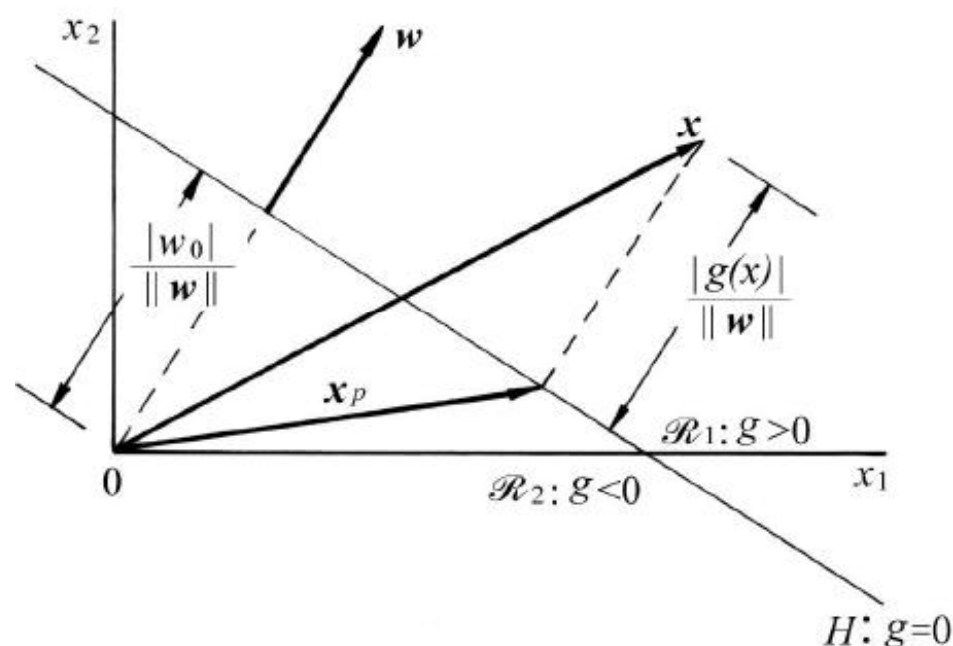
$$\text{即: } g(\mathbf{x}) = r \|\mathbf{w}\| \quad \text{或} \quad r = \frac{g(\mathbf{x})}{\|\mathbf{w}\|} \quad \|\mathbf{x} - \mathbf{x}_p\|_2 = \frac{|g(\mathbf{x})|}{\|\mathbf{w}\|}$$

即：判别函数  $g(\mathbf{x})$  正比于  $\mathbf{x}$  到超平面的代数距离  $r$

**问题：**如何确定任意观测样本到分类边界的距离？

$$g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + \omega_0 = r \|\mathbf{w}\|$$

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



对于原点  $\mathbf{x} = \mathbf{0}$ , 有:

$$g(\mathbf{x} = \mathbf{0}) = \omega_0, \quad r = \frac{\omega_0}{\|\mathbf{w}\|}$$

$$\begin{cases} \omega_0 > 0, & \text{原点 } \mathbf{x} = \mathbf{0} \text{ 在 } H \text{ 正侧;} \\ \omega_0 < 0, & \text{原点 } \mathbf{x} = \mathbf{0} \text{ 在 } H \text{ 负侧} \\ \omega_0 = 0, & H \text{ 经过原点, 判别函数 } g(\mathbf{x}) = \mathbf{w}^T \mathbf{x}. \end{cases}$$

**问题:** 如何确定  
原点到分类边界  
的距离?

超平面(分类面, 决策面) $H$ 的法向量为 $\boldsymbol{w}$ ;

超平面 $H$ 的偏移(超平面位置)由 $\omega_0$ 决定;

判别函数 $g(\boldsymbol{x})$ 正比于观测 $\boldsymbol{x}$ 到超平面 $H$ 的代数距离;

决策域 $\mathfrak{R}_1, \mathfrak{R}_2$ 分别位于超平面 $H$ 的正负侧。

**设计线性分类器的关键:**

基于何种准则, 估计判别函数的参数 $\boldsymbol{w}, \omega_0$ ;

# 主要内容

A. 引言

B. Fisher线性判别分析

C. 感知器

D. 线性支持向量机(线性SVM)

# B. Fisher判别分类模型

## 掌握：

- Fisher判别分类模型设计的基本思想
- 准则函数
- Fisher判别分类模型的实现步骤
- Fisher法进行监督式特征提取的前提条件

## 问题：

你能否将Fisher判别分类模型的“Fisher判别比”用于分类任务中的特征选择？

## 作业：

以鸢尾花数据集为例，采用Fisher判别分类模型(先监督式降维，再设计分类模型(如朴素贝叶斯))，进行分类。



# 主要内容

## 第一部分 类别数目 $C=2$ 时 FLDA 分类模型

关键：确定最佳投影直线

## 第二部分 类别数目 $C>2$ 时 FLDA 分类模型

许多机器学习问题涉及样本的“降维”

“降维”的目的不止一个，如：

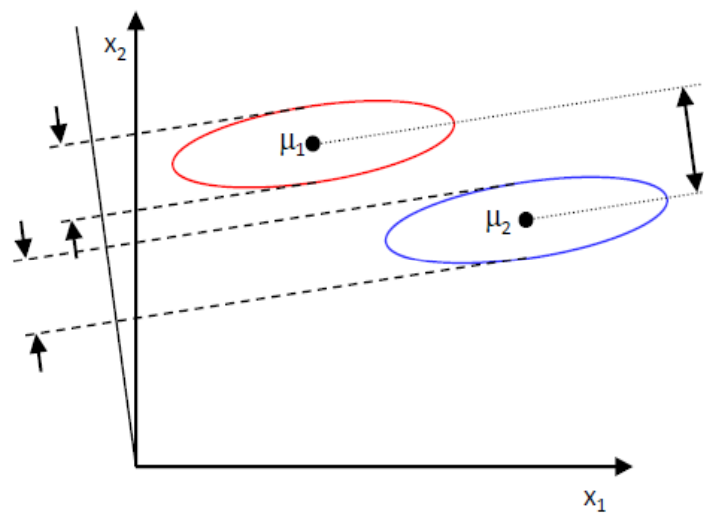
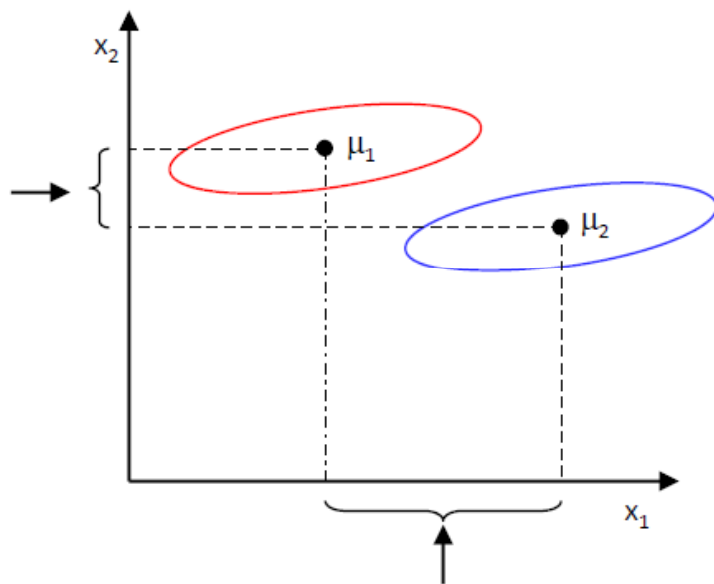
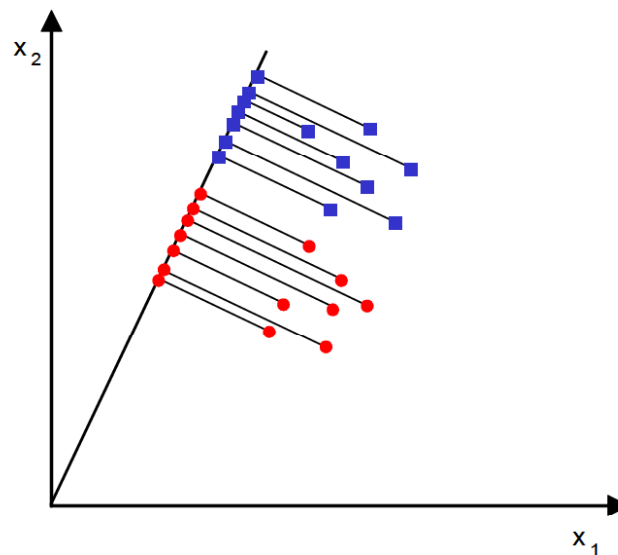
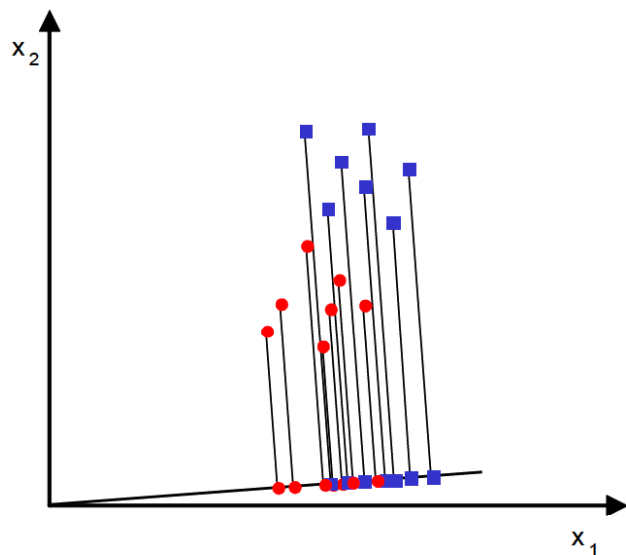
(1) *PCA*—**数据描述** (信息压缩)

目的：寻找有效表示特征的主轴方向

(2) *Fisher***线性判别**—**分类器设计**

目的：寻找能有效分类的方向

什么样的投影方向更有助于分类？



# 1 基本思想

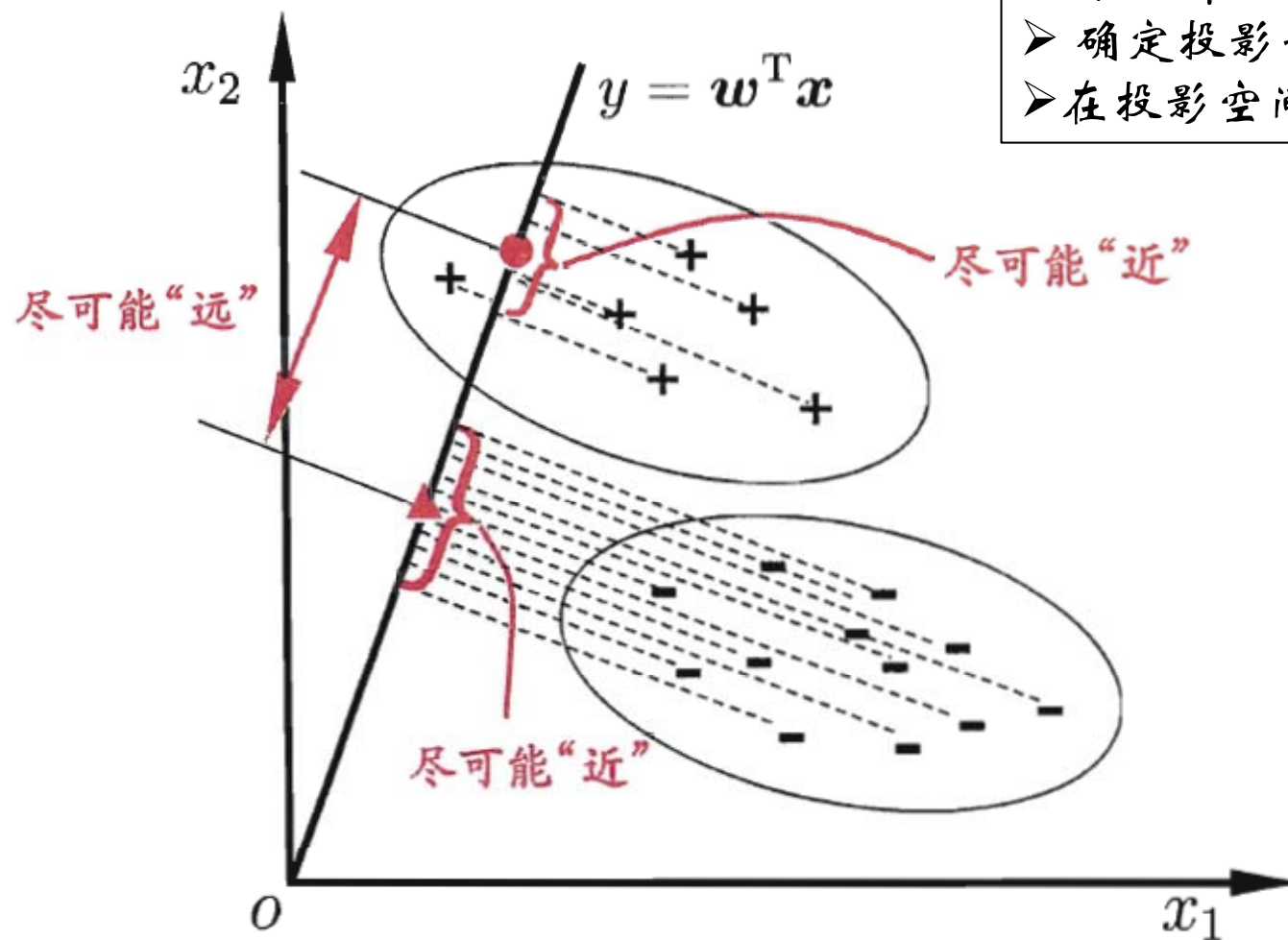
将 $d$ 维特征空间的样本向某低维子空间做投影  
(两类分类对应1维空间); 在该投影空间, 能最大限度区分各类数据点, 以有效分类.

{ 类间样本 尽可能远离  
  类内样本 尽可能聚集

分类器的设计  $\Rightarrow$  寻找投影直线的最佳方向 ( $C = 2$ )

{ 类间样本 尽可能远离  
类内样本 尽可能聚集

二维特征、两类别Fisher判别分析 示意  
引自周志华《机器学习》page60



任务分解:

- 确定投影子空间
- 在投影空间, 设计分类模型

## 2 问题描述

已知： $d$ 维训练样本  $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$

$$\text{其中} \begin{cases} \omega_1 \text{类: } \mathcal{X}_1 = \{\mathbf{x}_1^1, \dots, \mathbf{x}_{N_1}^1\} \\ \omega_2 \text{类: } \mathcal{X}_2 = \{\mathbf{x}_1^2, \dots, \mathbf{x}_{N_2}^2\} \end{cases}$$

$$\mathcal{X} = \mathcal{X}_1 \cup \mathcal{X}_2 \quad N = N_1 + N_2$$

$$\xrightarrow{\text{对 } \mathbf{x} \text{ 投影: } y = \mathbf{w}^T \mathbf{x}} \mathcal{Y} = \{y_1, \dots, y_N\} \quad \begin{cases} \omega_1 \text{类 } \mathcal{Y}_1, & N_1 \\ \omega_2 \text{类 } \mathcal{Y}_2, & N_2 \end{cases}$$

$$\text{其中 } y_i = \mathbf{w}^T \mathbf{x}_i, i = 1, 2, \dots, N$$

求解：最佳投影方向  $\mathbf{w}^*$ ，使两类分类效果最好。

### 3.准备工作 (定义几个必要的参量)

#### [1] 投影前 -- $d$ 维特征空间

#### 各类均值向量

$$m_i = \frac{1}{N_i} \sum_{x \in \mathcal{X}_i} x \quad i = 1, 2$$

#### 总体均值向量 (样本集的中心)

$$m = \frac{m_1 N_1 + m_2 N_2}{N_1 + N_2}$$

#### 类内散度矩阵 (*within-class scatter matrix*)

$$S_i = \frac{1}{N_i} \sum_{x \in \mathcal{X}_i} (x - m_i)(x - m_i)^T \quad i = 1, 2$$

## [1] 投影前 -- $d$ 维特征空间

总类内散度矩阵 (*pooled within-class scatter matrix*)

$$S_w = \frac{1}{N_1 + N_2} \sum_{i=1}^2 N_i S_i$$

$S_w$  对称半正定；若  $N > d$ ,  $S_w$  通常非奇异

类间散度矩阵 (*between-class scatter*)

$$\begin{cases} \text{不考虑先验概率} & S_b = (m_1 - m_2)(m_1 - m_2)^T \\ \text{若考虑先验概率} & S_b = P(\omega_1)P(\omega_2)(m_1 - m_2)(m_1 - m_2)^T \end{cases}$$

$$S_b = \frac{1}{N_1 + N_2} \sum_{i=1}^2 N_i (m_i - m)(m_i - m)^T = \sum_{i=1}^2 P(\omega_i) (m_i - m)(m_i - m)^T$$

或者 
$$S_b = \frac{1}{2} \sum_{i=1}^2 (m_i - m)(m_i - m)^T$$



## [2]投影后-1维空间

### 样本均值

$$\widetilde{m}_i = \frac{1}{N_i} \sum_{y \in \mathcal{Y}_i} y = \frac{1}{N_i} \sum_{x \in \mathcal{X}_i} w^T x = w^T \left( \frac{1}{N_i} \sum_{x \in \mathcal{X}_i} x \right) = w^T m_i$$
$$i = 1, 2$$

### 总体样本均值

$$\widetilde{m} = \frac{\widetilde{m}_1 N_1 + \widetilde{m}_2 N_2}{N_1 + N_2} = w^T m$$

### 类内离散度

$$\begin{aligned} \widetilde{S}_i^2 &= \frac{1}{N_i} \sum_{y_j \in \mathcal{Y}_i} (y_j - \widetilde{m}_i)^2 \\ &= \frac{1}{N_i} \sum_{x_j \in \mathcal{X}_i} w^T (x_j - m_i) (x_j - m_i)^T w = w^T S_i w \quad i = 1, 2 \end{aligned}$$

## [2]投影后-1维空间

### 总类内离散度

$$\tilde{S}_w = \tilde{S}_1^2 + \tilde{S}_2^2 = w^T S_1 w + w^T S_2 w = w^T (S_1 + S_2) w = w^T S_w w$$

$$\tilde{S}_w = \frac{1}{N_1 + N_2} \sum_{i=1}^2 N_i \tilde{S}_i^2 = \frac{1}{N_1 + N_2} \sum_{i=1}^2 N_i w^T S_i w = w^T S_w w$$

### 类间离散度

$$\begin{aligned} \tilde{S}_b^2 &= (\tilde{m}_1 - \tilde{m}_2)^2 = (w^T m_1 - w^T m_2)^2 \\ &= (w^T m_1 - w^T m_2)(m_1^T w - m_2^T w) \\ &= w^T (m_1 - m_2)(m_1 - m_2)^T w = w^T S_b w \end{aligned}$$

$$\begin{aligned} \tilde{S}_b^2 &= \frac{1}{N_1 + N_2} \sum_{i=1}^2 N_i (\tilde{m}_i - \tilde{m})(\tilde{m}_i - \tilde{m})^T \\ &= \frac{1}{N_1 + N_2} \sum_{i=1}^2 N_i w^T (m_i - m)(m_i - m)^T w = w^T S_b w \end{aligned}$$

## 4.构造Fisher准则函数

高维( $d$ 维)样本投影到1维空间

{ 类间样本尽量分开  $\Rightarrow$  投影后类间离散度  $\uparrow$   
类内样本尽量聚集  $\Rightarrow$  投影后总类内离散度  $\downarrow$

$$\Rightarrow \begin{array}{l} \text{Fisher准则函数} \quad J_F(\mathbf{w}) = \frac{|\tilde{m}_1 - \tilde{m}_2|^2}{\tilde{S}_w} = \frac{(\tilde{m}_1 - \tilde{m}_2)^2}{\tilde{S}_w} \\ \text{Fisher准则} \quad \max_{\mathbf{w}} J_F(\mathbf{w}) \\ \mathbf{w}^* = \arg \max_{\mathbf{w}} J_F(\mathbf{w}) \quad \text{如何求解?} \end{array}$$

Fisher线性分类器设计 { (1) 最佳投影空间的估计  
(2) 基于训练样本在最佳投影空间  
投影设计最优边界。

## 5.确定最优投影方向--求解 $\mathbf{w}^* = \arg \max_{\mathbf{w}} J_F(\mathbf{w})$

### [1]确定 $J_F(\mathbf{w})$ 关于 $\mathbf{w}$ 的显式函数形式

$$J_F(\mathbf{w}) = \frac{(\tilde{m}_1 - \tilde{m}_2)^2}{\tilde{S}_1^2 + \tilde{S}_2^2}$$

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

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

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

所以  $\boxed{J_F(\mathbf{w}) = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_w \mathbf{w}}}$  ----- 广义瑞利(Rayleigh) 商.

$$[2] \text{ 估计 } \mathbf{w}^* = \arg \max_{\mathbf{w}} \mathbf{J}_F(\mathbf{w}) = \arg \max_{\mathbf{w}} \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_w \mathbf{w}}$$

$$\Rightarrow \boxed{\begin{array}{ll} \max_{\mathbf{w}} & \mathbf{w}^T \mathbf{S}_B \mathbf{w} \\ \text{s.t.} & \mathbf{w}^T \mathbf{S}_w \mathbf{w} = c \neq 0 \end{array}}$$

采用**Lagrange**条件极值法，引入辅助函数：

$$\begin{aligned} L(\mathbf{w}, \lambda) &= \mathbf{J}_F(\mathbf{w}) - \lambda(\mathbf{w}^T \mathbf{S}_w \mathbf{w} - c) \\ &= \mathbf{w}^T \mathbf{S}_B \mathbf{w} - \lambda(\mathbf{w}^T \mathbf{S}_w \mathbf{w} - c) \end{aligned}$$

$\lambda$ 为标量**Lagrange**乘子

$$L(\mathbf{w}, \lambda) = \mathbf{w}^T \mathbf{S}_B \mathbf{w} - \lambda (\mathbf{w}^T \mathbf{S}_w \mathbf{w} - c)$$

$$\frac{\partial L}{\partial \mathbf{w}} = 2(\mathbf{S}_B \mathbf{w} - \lambda \mathbf{S}_w \mathbf{w})$$

$$\text{令 } \frac{\partial L}{\partial \mathbf{w}} = 0$$

$$\text{得 } \mathbf{S}_B \mathbf{w}^* = \lambda \mathbf{S}_w \mathbf{w}^* \xrightarrow{N > d, \mathbf{S}_w \text{非奇异}} \boxed{\mathbf{S}_w^{-1} \mathbf{S}_B \mathbf{w}^* = \lambda \mathbf{w}^*}$$

$$\text{相应地: } J(\mathbf{w}^*) = \frac{\mathbf{w}^{*T} \mathbf{S}_B \mathbf{w}^*}{\mathbf{w}^{*T} \mathbf{S}_w \mathbf{w}^*} = \frac{\lambda \mathbf{w}^{*T} \mathbf{S}_w \mathbf{w}^*}{\mathbf{w}^{*T} \mathbf{S}_w \mathbf{w}^*} = \lambda$$

**思路1** 要使  $J(\mathbf{w}^*) = \lambda$  最大,

应使  $\lambda$  取  $\mathbf{S}_w^{-1} \mathbf{S}_B$  矩阵的最大本征值,

$\mathbf{w}^*$  为  $\mathbf{S}_w^{-1} \mathbf{S}_B$  矩阵相应本征列向量

$$\frac{\partial (\mathbf{x}^T \mathbf{A} \mathbf{x})}{\partial \mathbf{x}} = (\mathbf{A} + \mathbf{A}^T) \mathbf{x}$$

$\mathbf{A}$  是对称矩阵  
 $\quad \quad \quad = 2\mathbf{A}\mathbf{x}$

[2] **确定**  $\boldsymbol{w}^* = \arg \max_{\boldsymbol{w}} J_F(\boldsymbol{w})$

**思路2: 直接找到**  $\boldsymbol{w}^*$

类间散度矩阵  $\boldsymbol{S}_B = (\boldsymbol{m}_1 - \boldsymbol{m}_2)(\boldsymbol{m}_1 - \boldsymbol{m}_2)^T$

$$\boldsymbol{S}_B \boldsymbol{w}^* = (\boldsymbol{m}_1 - \boldsymbol{m}_2) \underbrace{(\boldsymbol{m}_1 - \boldsymbol{m}_2)^T \boldsymbol{w}^*}_{\text{标量 } R} = (\boldsymbol{m}_1 - \boldsymbol{m}_2) R$$

所以  $\boldsymbol{S}_B \boldsymbol{w}^*$  方向与向量  $(\boldsymbol{m}_1 - \boldsymbol{m}_2)$  方向一致。

$$\lambda \boldsymbol{w}^* = \boldsymbol{S}_w^{-1} \boldsymbol{S}_B \boldsymbol{w}^* = \boldsymbol{S}_w^{-1} (\boldsymbol{m}_1 - \boldsymbol{m}_2) R$$

$$\Rightarrow \boldsymbol{w}^* = \frac{R}{\lambda} \boldsymbol{S}_w^{-1} (\boldsymbol{m}_1 - \boldsymbol{m}_2)$$

忽略标量  $\frac{R}{\lambda}$ , 得 **最优投影方向**  $\boldsymbol{w}^* = \boldsymbol{S}_w^{-1} (\boldsymbol{m}_1 - \boldsymbol{m}_2)$

### [3] $w^* = \arg \max_w J_F(w)$ 计算流程

**STEP1.** 获取来自两类  $\omega_1 / \omega_2$  的训练样本集  $\mathcal{X}$

$$\begin{cases} \omega_1 \text{类: } \mathcal{X}_1, N_1 \text{个样本} \\ \omega_2 \text{类: } \mathcal{X}_2, N_2 \text{个样本} \end{cases}$$

**STEP2.** 计算各类样本均值  $m_i I = \frac{1}{N_i} \sum_{x \in \mathcal{X}_i} x \quad i = 1, 2$

**STEP3.** 类内散度矩阵

$$S_i = \frac{1}{N_i} \sum_{x \in \mathcal{X}_i} (x - m_i)(x - m_i)^T \quad i = 1, 2$$

**STEP4.** 总类内散布矩阵  $S_w$  及逆  $S_w^{-1}$ .

$$S_w = \frac{1}{N_1 + N_2} \sum_{i=1}^2 N_i S_i$$

为避免  $S_w$  的奇异,  
引入一个小正数  $\varepsilon$

计算  $w^* = (S_w + \varepsilon I)^{-1} (m_1 - m_2)$

**STEP5.** 计算  $w^*$   $w^* = S_w^{-1} (m_1 - m_2)$

利用上述过程, 可直接确定投影直线  $w^*$ .



## 6.确定最优分类超平面

$$(\boldsymbol{w}^*)^T \boldsymbol{x} + \omega_0 = 0$$

明确  $\begin{cases} \text{法向量} \\ \omega_0 = ? \end{cases} \quad \boldsymbol{w}^* = \arg \max_{\boldsymbol{w}} J_F(\boldsymbol{w})$

# 一维最佳投影空间的分类边界

## 基于最小错误率贝叶斯决策

第1， 两类样本正态分布， 且分布形状一致  
——线性分类器

第2， 若维数 $d$ 及样本数目 $N$ 足够大时，  $y = \mathbf{w}^{*T} \mathbf{x}$   
近似正态分布， 可在投影空间利用“**两步  
贝叶斯决策**”法。

## A. 确定阈值点(分类边界) $y_0$

先验概率相等  $y_0 = \frac{\tilde{m}_1 + \tilde{m}_2}{2} = \frac{(m_1 + m_2) \mathbf{S}_w^{-1} (m_1 - m_2)}{2}$

先验概率不等

$$y_0 = \frac{\tilde{m}_1 + \tilde{m}_2}{2} + \log \frac{P(\omega_2)}{P(\omega_1)} = \frac{(m_1 + m_2) \mathbf{S}_w^{-1} (m_1 - m_2)}{2} + \log \frac{P(\omega_2)}{P(\omega_1)}$$

投影后样本数据中心

$$y_0 = \frac{\tilde{m}_1 N_1 + \tilde{m}_2 N_2}{N_1 + N_2} = \frac{(m_1 N_1 + m_2 N_2) \mathbf{S}_w^{-1} (m_1 - m_2)}{N_1 + N_2}$$

## B.决策

对于任意观测 $\mathbf{x}$

$$\left\{ \begin{array}{l} \text{投影空间: } y = \mathbf{w}^{*T} \mathbf{x} \\ \text{判决: 若 } y > y_0, \text{ 则 } \mathbf{x} \text{ 为 } \omega_1 \text{ 类;} \\ \quad \text{若 } y < y_0, \text{ 则 } \mathbf{x} \text{ 为 } \omega_2 \text{ 类。} \end{array} \right.$$

分类边界满足:  $y - y_0 = 0$ , 即  $\mathbf{w}^{*T} \mathbf{x} - y_0 = 0$

### 原始特征空间

判别函数:  $g(\mathbf{x}) = \mathbf{w}^{*T} \mathbf{x} + \omega_0 = \mathbf{w}^{*T} \mathbf{x} - y_0$

分类边界:  $\mathbf{w}^{*T} \mathbf{x} - y_0 = 0$

决策规则:  $\left\{ \begin{array}{l} \text{若 } g(\mathbf{x}) > 0, \text{ 则 } \mathbf{x} \text{ 为 } \omega_1 \text{ 类} \\ \text{若 } g(\mathbf{x}) < 0, \text{ 则 } \mathbf{x} \text{ 为 } \omega_2 \text{ 类} \end{array} \right.$

# 主要内容

## 第一部分 类别数目 $C=2$ 时 FLDA 分类模型

关键: 确定最佳投影直线

## 第二部分 类别数目 $C>2$ 时 FLDA 分类模型

关键: 确定最佳投影矩阵

注意: 线性投影子空间的维数最高为  $C-1$

# 1 问题描述

已知： $d$ 维训练样本  $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\} = \mathcal{X}_1 \cup \mathcal{X}_2 \cup \dots \cup \mathcal{X}_C$

其中  $\omega_i$ 类,  $\mathcal{X}_i$ ,  $N_i$ 个

$$N = \sum_{i=1}^C N_i$$

$$\xrightarrow{\text{对}\mathbf{x}\text{投影: } y=W^T x} \mathcal{Y} = \{\mathbf{y}_1, \dots, \mathbf{y}_N\} \quad \left\{ \begin{array}{l} \omega_1 \text{类 } \mathcal{Y}_1, N_1 \\ \vdots \\ \omega_C \text{类 } \mathcal{Y}_C, N_C \end{array} \right.$$

其中  $\mathbf{y}_j = \mathbf{W}^T \mathbf{x}_j, i = 1, 2, \dots, N$

求解：最佳投影子空间  $\mathbf{W}^*$ ，使 $C$ 类分类效果最好。

## 2.准备工作

### [1]投影前-- $d$ 维特征空间

各类均值向量 
$$\mathbf{m}_i = \frac{1}{N_i} \sum_{x \in \mathcal{X}_i} \mathbf{x} \quad i = 1, \dots, C$$

总体均值向量 
$$\mathbf{m} = \sum_{i=1}^C \frac{N_i}{N} \mathbf{m}_i$$

类内散度矩阵 (*within-class scatter matrix*)

$$\mathbf{S}_i = \frac{1}{N_i} \sum_{x \in \mathcal{X}_i} (\mathbf{x} - \mathbf{m}_i)(\mathbf{x} - \mathbf{m}_i)^T \quad i = 1, \dots, C$$

总类内散度矩阵 (*pooled within-class scatter matrix*)

$$\mathbf{S}_w = \sum_{i=1}^C \frac{N_i}{N} \mathbf{S}_i$$

$\mathbf{S}_w$  对称半正定；若  $N > d$ ,  $\mathbf{S}_w$  通常非奇异

类间散度矩阵 (*between-class scatter*) 
$$\mathbf{S}_b = \sum_{i=1}^C \frac{N_i}{N} (\mathbf{m}_i - \mathbf{m})(\mathbf{m}_i - \mathbf{m})^T$$

## [2] 投影后 – 低维空间 (维数 $\leq C - 1$ )

总类内散度矩阵  $\tilde{S}_w = W^T S_w W$

$$\tilde{S}_w = W^T S_w W = W^T \left( \sum_{i=1}^C \frac{N_i}{N} S_i \right) W$$

类间散度矩阵  $\tilde{S}_b = W^T S_b W$

$$\tilde{S}_b = W^T S_b W = W^T \left[ \sum_{i=1}^C \frac{N_i}{N} (m_i - m)(m_i - m)^T \right] W$$



### 3. Fisher 准则函数

Fisher 准则函数  $J_F(W) = \frac{\det(\tilde{S}_b)}{\det(\tilde{S}_w)} = \frac{\det(W^T S_b W)}{\det(W^T S_w W)}$

或者  $J_F(W) = \frac{\text{tr}(\tilde{S}_b)}{\text{tr}(\tilde{S}_w)} = \frac{\text{tr}(W^T S_b W)}{\text{tr}(W^T S_w W)}$

类间远离  
类内紧致

或者  $J_F(W) = \frac{\prod_{diag} \tilde{S}_b}{\prod_{diag} \tilde{S}_w} = \frac{\prod_{diag} W^T S_b W}{\prod_{diag} W^T S_w W} = \frac{\prod_j w_j^T S_b w_j}{\prod_j w_j^T S_w w_j} = \prod_j \frac{w_j^T S_b w_j}{w_j^T S_w w_j}$

$S_b$  的秩最大值为  $C-1$

$\prod_{diag} A$  — 矩阵  $A$  的主对角线元素乘积

Fisher 准则  $\max_{W \in R^{d \times (C-1)}} J_F(W)$

最优投影子空间  $W^* = \arg \max_{W \in R^{d \times (C-1)}} J_F(W)$

并且  $W^* = [w_1 \cdots w_{C-1}]_{d \times (C-1)}$

#### 4.确定最优投影子空间--求解 $W^* = \arg \max_W J_F(W)$

$$\text{例: } J_F(W) = \frac{\text{tr}(\tilde{S}_b)}{\text{tr}(\tilde{S}_w)} = \frac{\text{tr}(W^T S_b W)}{\text{tr}(W^T S_w W)}$$

$$\max_W J_F(W) \Rightarrow \begin{cases} \max_W \text{tr}(W^T S_b W) \\ \text{s.t. } \text{tr}(W^T S_w W) = 1 \neq 0 \end{cases}$$

$$\Rightarrow \begin{cases} L(W, \lambda) = \text{tr}(W^T S_b W) - \lambda(\text{tr}(W^T S_w W) - 1) \\ \max_{W, \lambda} L(W, \lambda) \end{cases}$$

$$L(W, \lambda) = \text{tr}(W^T S_b W) - \lambda(\text{tr}(W^T S_w W) - 1)$$

$$\frac{\partial L(W, \lambda)}{\partial W} = 2(S_b W - \lambda S_w W) = 0$$

$$\frac{\partial \text{tr}(X^T A X)}{\partial X} = [X^T (A^T + A)]^T = (A + A^T) X$$

$$S_b W = \lambda S_w W = S_w W \Lambda \quad S_b w_i = \lambda S_w w_i \quad i = 1, \dots, C-1$$

$$\mathbf{S}_b \mathbf{W} = \lambda \mathbf{S}_w \mathbf{W}$$

$$\mathbf{S}_b \mathbf{w}_i = \lambda \mathbf{S}_w \mathbf{w}_i \quad i = 1, \dots, C-1$$

$$\mathbf{S}_w^{-1} \mathbf{S}_b \mathbf{w}_i = \lambda \mathbf{w}_i$$

为避免 $\mathbf{S}_w$ 的奇异，  
引入一个小正数 $\varepsilon$   
计算 $(\mathbf{S}_w + \varepsilon \mathbf{I})^{-1} \mathbf{S}_b$

若 $\lambda$ 为矩阵 $\mathbf{S}_w^{-1} \mathbf{S}_b$ 的某个本征值，则 $\mathbf{w}_i$ 就是相应的本征列向量。

若有 $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_{C-1} > 0$

相应 $\mathbf{w}_1^*, \mathbf{w}_2^*, \dots, \mathbf{w}_{C-1}^*$ 构成最佳投影矩阵

$$\mathbf{W}^* = \begin{bmatrix} \mathbf{w}_1^* & \mathbf{w}_2^* & \dots & \mathbf{w}_{C-1}^* \end{bmatrix}$$

这就是 $C-1$ 维投影子空间。

任意观测样本 $\mathbf{x} \in \mathfrak{R}$ ，特征提取结果： $\mathbf{y} = \mathbf{W}^{*T} \mathbf{x}$

# 主要内容

A. 引言

B. Fisher线性判别分析

C. 感知器

D. 线性支持向量机(线性SVM)

## B. 感知器分类模型

要求:

- 理解感知器准则函数的构造思想
- 掌握感知器准则函数的形式、意义、
- 感知器分类模型的学习
- 基于感知器的类别决策
- 感知器分类模型学习的前提条件

关键词:

- 解向量、解空间
- 训练样本的增广、规范化
- 权向量的增广
- 线性可分的训练样本
- 梯度下降法(固定步长/变步长; 单样本/批量样本)

## 1. 基本思想

针对两类别线性可分问题，引入解向量；

基于错分的训练样本，构造感知器准则函数；

采用梯度下降法，估计解向量，得线性分类模型。

## 2.几个基本概念

### (1)增广的特征向量 / 权向量

线性判别函数  $g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + \omega_0 = \mathbf{w}^T \mathbf{x} + \mathbf{1}\omega_0$

定义  $\left\{ \begin{array}{l} \text{增广的特征向量 } \mathbf{y} = \begin{bmatrix} 1 & x_1 & \cdots & x_d \end{bmatrix}^T = \begin{bmatrix} \mathbf{1} \\ \mathbf{x} \end{bmatrix} \\ \text{增广的权向量 } \mathbf{a} = \begin{bmatrix} \omega_0 & w_1 & \cdots & w_d \end{bmatrix}^T = \begin{bmatrix} \omega_0 \\ \mathbf{w} \end{bmatrix} \end{array} \right.$

→  $\left\{ \begin{array}{l} \text{线性判别函数 } g(\mathbf{y}) = \mathbf{a}^T \mathbf{y} \\ \text{决策规则} \quad \text{若 } g(\mathbf{y}) > 0, \text{ 则决策 } \mathbf{y} \in \omega_1 \\ \quad \quad \quad \text{若 } g(\mathbf{y}) < 0, \text{ 则决策 } \mathbf{y} \in \omega_2 \end{array} \right.$

## (2) 训练样本集的线性可分性

考虑  $\begin{cases} \text{两类: } \omega_1, \omega_2 \\ \text{增广样本集 } \mathcal{Y} = \{\mathbf{y}_1, \dots, \mathbf{y}_N\} \end{cases}$

若至少存在一个线性分类器  $g(\mathbf{y}) = \mathbf{a}^T \mathbf{y}$ , 能使所有训练样本正确分类, 即

$$\exists \mathbf{a}, \begin{cases} \text{若 } \mathbf{y}_i \in \omega_1 \text{ 类, 则有 } \mathbf{a}^T \mathbf{y}_i > 0 \\ \text{若 } \mathbf{y}_i \in \omega_2 \text{ 类, 则有 } \mathbf{a}^T \mathbf{y}_i < 0 \end{cases} \quad i = 1, 2, \dots, N$$

则上述训练样本是**线性可分**的; 否则, **线性不可分**。

若样本集  $\mathcal{Y}$  线性可分, 则必存在向量  $\mathbf{a}$ , 对各训练样本正确分类。



### (3) 规范化增广样本向量

对于**线性可分的增广样本集** $\mathcal{Y} = \{y_1, \dots, y_N\}$ ,

定义**规范化增广样本向量** $y'$ , 满足

$$y_i' = \begin{cases} y_i & \text{若 } y_i \in \omega_1 \text{ 类} \\ -y_i & \text{若 } y_i \in \omega_2 \text{ 类} \end{cases} \quad i = 1, \dots, N$$

则必存在权向量 $a$ , 使得

$$a^T y_i' > 0, \quad i = 1, \dots, N$$

目的在于: 模型学习过程中, 有助于区分正确分类的训练样本与错分的训练样本

可记**规范化增广样本向量** $y'$ 为 $y$

## (4) 解向量和解区

对于线性可分的训练样本集，设其**规范化增广样本集**为  $\mathcal{Y} = \{y_1, \dots, y_N\}$ ，若存在权向量  $a$  满足

$$a^T y_i > 0, \quad i = 1, \dots, N$$

则称权向量  $a$  为1个**解向量**，记为  $a^*$ 。

**权值空间**中，所有解向量组成的区域称为**解区**。

## (4) 解向量和解区 -- 续

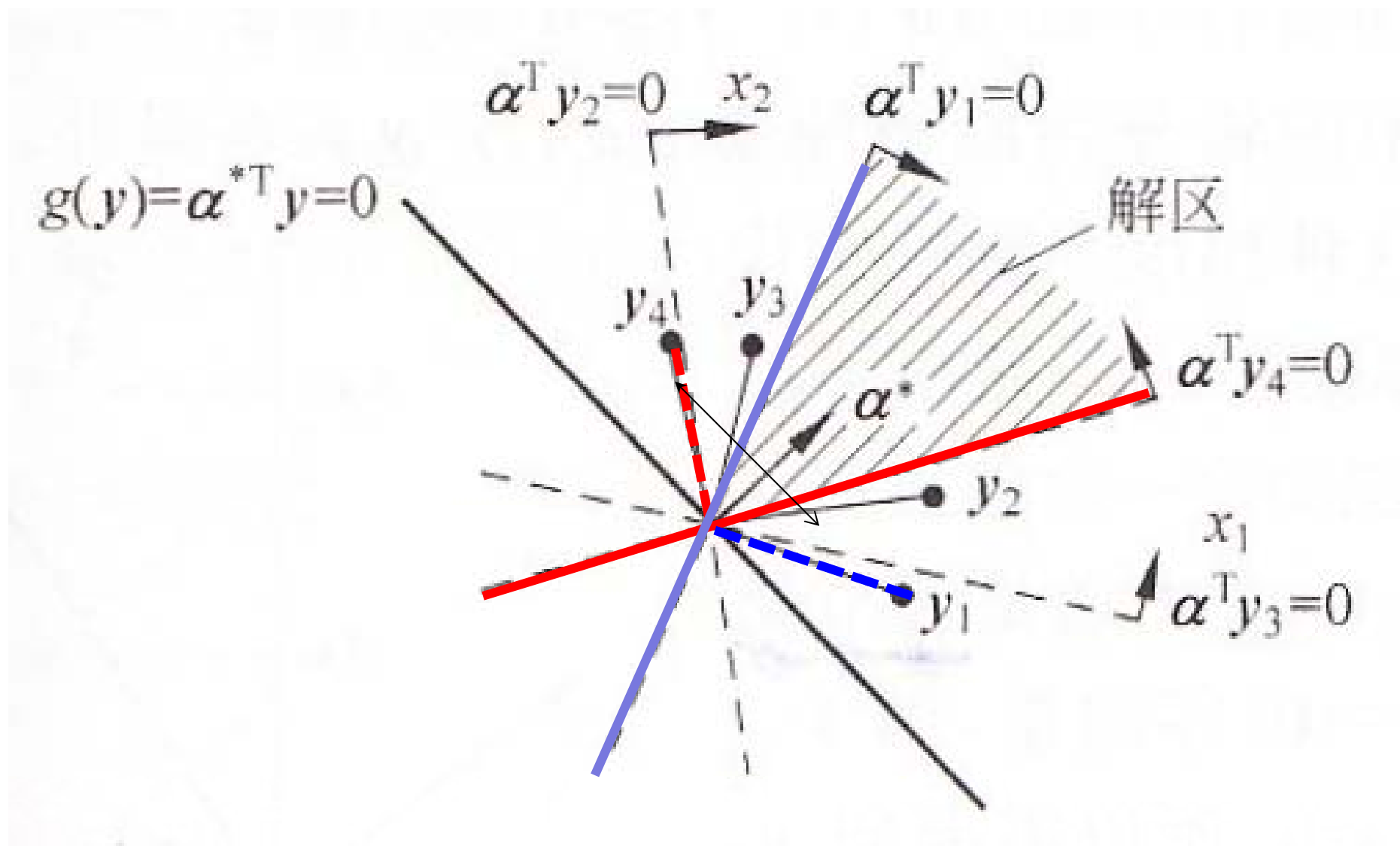
对于每个样本  $y_i$

$$\left\{ \begin{array}{l} a^T y_i = 0 \text{ 是在权值空间过原点的超平面 } \hat{H}_i, \\ \text{超平面 } \hat{H}_i \text{ 的法向量是 } y_i; \\ \text{解向量 } a^* \text{ 在超平面正侧 (权值空间正侧 } a^T y_i > 0) \end{array} \right.$$

$N$  个样本产生  $N$  个超平面

**解区** 是 **权值空间**  $N$  个超平面 **正侧** 区域的交集。

## 规范化增广样本集、解向量和解区



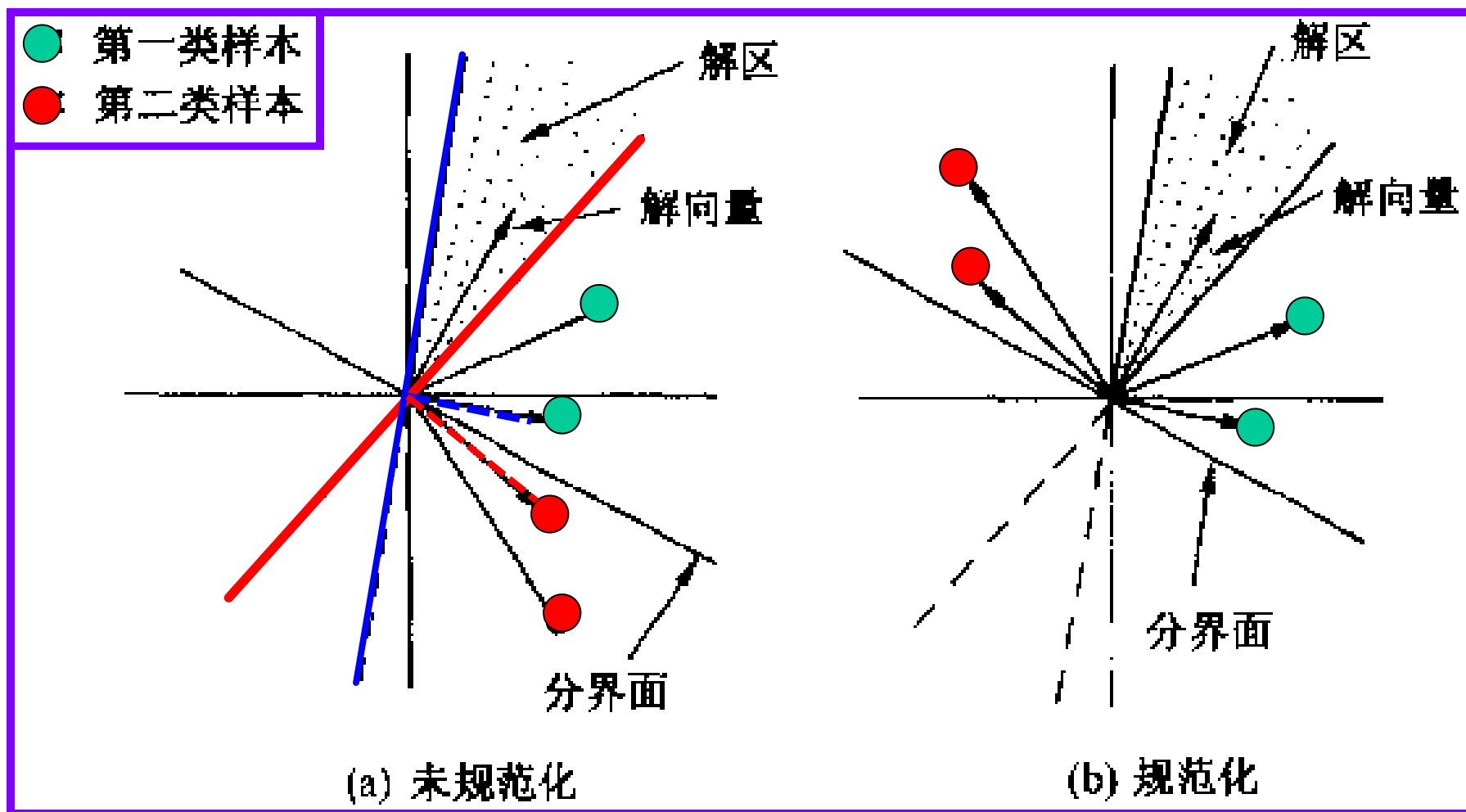


图 二维空间、两类情况下，增广训练样本规范化前后，训练样本及其特征空间(与权空间维数一致)的解区域及解向量示意图。

## (5) 对解区的限制 / 余量

若解向量存在，通常不是惟一的。

解区所有的解向量都满足  $a^T y_i > 0$ .

通常，解向量越靠近解区中间

→ 分类面离各类样本越远

→ 它到各样本的最近距离越远

→ 对新样本错分的可能性越小

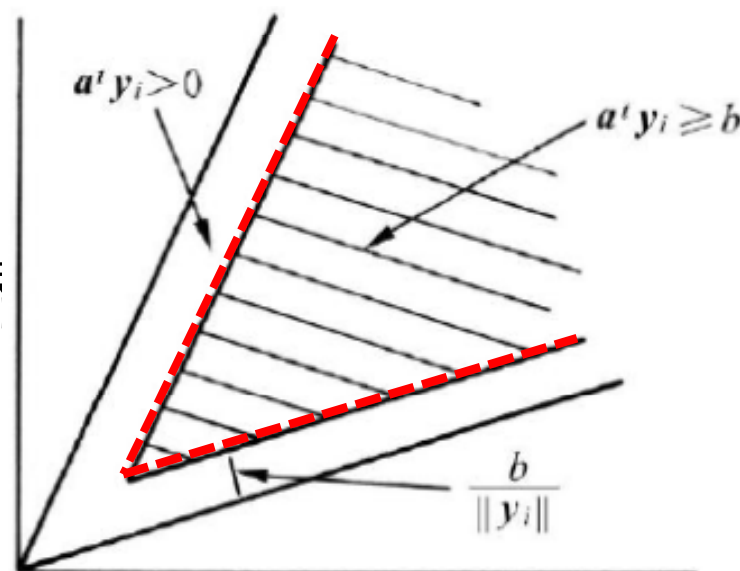
⇒ 为使解向量更可靠，需要**限制解区**。

收缩解区，只考虑**靠近解区中间**的解向量。

引入“**余量**(*margin*)” $b > 0$ , 在解区中寻找解向量 $\mathbf{a}^*$ ,

$$\text{满足 } \mathbf{a}^T \mathbf{y}_i > b > 0, \quad i = 1, \dots, N$$

好处 { 解向量更可靠  
推广性更强  
防止算法收敛至解区边界



**引入余量的解区**

### 3.感知器准则函数的构建 / 解向量求解

#### (1)构建感知器准则函数 $J_p(a)$

给定**规范化增广**样本集

$$\mathcal{Y} = \{y_1, \dots, y_N\}$$

若权向量 $a$ 使得训练样本 $y_k$ 被错分, 则有  $a^T y_k < 0$

对于**所有错分样本**, 定义**感知器准则函数**

$$J_p(a) = \sum_{a^T y_k \leq 0} (-a^T y_k)$$

函数值越大, 训练集的  
错分程度越高

当且仅当 $J_p(a^*) = \min J_p(a) = 0$ , 错分样本数目为0,

$a^*$ 为解向量.

线性分类模型的学习目的, 就是如何利用  
规范增广的训练集, 估计这个**解向量**



(2) 确定解向量  $\mathbf{a}^* = \arg \min_{\mathbf{a}} J_p(\mathbf{a}) = \sum_{\mathbf{a}^T \mathbf{y}_k \leq 0} (-\mathbf{a}^T \mathbf{y}_k)$

采用梯度下降法 (*gradient decent method*)

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

迭代形式  $\mathbf{a}(t+1) = \mathbf{a}(t) - \rho_t \cdot \nabla J_p(\mathbf{a}(t))$   
 $= \mathbf{a}(t) + \rho_t \cdot \sum_{\mathbf{a}^T(t) \cdot \mathbf{y}_k \leq 0} \mathbf{y}_k$

$\mathbf{a}(t+1)$  处准则函数取值:

$$J_p(\mathbf{a}(t+1)) = \sum_{\mathbf{a}^T(t+1) \cdot \mathbf{y}_k \leq 0} [-\mathbf{a}^T(t+1) \cdot \mathbf{y}_k]$$

由  $\mathbf{a}(0)$  开始得权向量序列:  $\mathbf{a}(1), \mathbf{a}(2), \dots, \mathbf{a}(t), \mathbf{a}(t+1), \dots$

## 方法1：固定增量单样本法

(*Fixed – Increment Single Sample Method*)

将样本集  $\{y_1, \dots, y_N\}$  视为一个不断重复出现的样本序列, 逐样本考察。

学习率(决定步长)固定为常数  $\rho_t = \rho (=1) > 0$

$$a(t+1) = a(t) - \rho \cdot \nabla J_p(a(t)) = a(t) + y_k$$

**STEP1.** 任意选择初始权向量  $a(0)$ ,  $t \leftarrow 0$ ;  $\rho (=1) > 0$

**STEP2.** 顺次考察各样本.

若  $[a(t)]^T y_j \leq 0$ , 则修正  $a(t)$  为  $a(t+1) = a(t) + y_j$

**STEP3.** 重复STEP2, 直到  $J_p(a) = 0$ , 得解向量  $a^*$

$$\begin{aligned}\left[\boldsymbol{a}(t+1)\right]^T \boldsymbol{y}_k &= \left[\boldsymbol{a}(t) + \boldsymbol{y}_k\right]^T \boldsymbol{y}_k \\ &= \left[\boldsymbol{a}(t)\right]^T \boldsymbol{y}_k + \left\|\boldsymbol{y}_k\right\|^2 \geq \left[\boldsymbol{a}(t)\right]^T \boldsymbol{y}_k\end{aligned}$$

**只要解区存在，逐渐修正权向量，不断接近超平面正侧**

**对于线性可分样本集，经有限步迭代，总可找到解向量 $\boldsymbol{a}^*$ 。**

## 方法2-带余量的可变增量单样本修正算法

考虑余量  $b > 0$   $\begin{cases} \text{严格错分条件} & [a(t)]^T y_k \leq b \\ \text{学习率} & \rho_t > 0 \text{ 可变 (如何自适应调整 } \rho_t?) \end{cases}$

对于  $a(t)$ , 若  $[a(t)]^T y_k \leq b$ , 则修正  $a(t)$

得  $a(t+1) = a(t) + \rho_t y_k$

要使  $[a(t+1)]^T y_k > b > 0$

$$\begin{aligned} \text{则 } [a(t+1)]^T y_k &= [a(t) + \rho_t y_k]^T y_k \\ &= [a(t)]^T y_k + \rho_t \cdot \|y_k\|^2 > b > 0 \end{aligned}$$

$$\Rightarrow \text{可变步长 } \rho_t > \frac{b - [a(t)]^T y_k}{\|y_k\|^2}$$

思考：若余量为0，  
如何自适应调整  $\rho_t$ ？

## 4.基于感知器的样本分类实现步骤

$$\text{判别函数} \begin{cases} g(\mathbf{x}; \mathbf{w}, \omega_0) = \mathbf{w}^T \mathbf{x} + \omega_0 = \mathbf{w}^T \mathbf{x} + 1\omega_0 \\ g(\mathbf{y}; \mathbf{a}) = \mathbf{a}^T \mathbf{y} = \begin{bmatrix} \omega_0 & \mathbf{w}^T \end{bmatrix} \begin{bmatrix} 1 \\ \mathbf{x} \end{bmatrix} \end{cases}$$

若训练集不是线性可分，  
将不能找到解向量

**STEP1.**原始特征空间训练样本集的**增广、规范化**。

**STEP2.**分类模型的学习。

利用**增广、规范化**训练集估计**解向量** $\mathbf{a}^* = \begin{bmatrix} \omega_0^* \\ \mathbf{w}^* \end{bmatrix}$

得 $\mathbf{w}^*$ ， $\omega_0^*$ ，判别函数为： $g(\mathbf{x}) = \mathbf{w}^{*T} \mathbf{x} + \omega_0^*$

**STEP3.**分类模型的使用。

对于原始特征空间的任意观测样本 $\mathbf{x}$ ，决策规则

$\begin{cases} \text{若 } g(\mathbf{x}) = \mathbf{w}^{*T} \mathbf{x} + \omega_0^* > 0, \text{ 则判断 } \mathbf{x} \text{ 为第1类} \\ \text{若 } g(\mathbf{x}) = \mathbf{w}^{*T} \mathbf{x} + \omega_0^* < 0, \text{ 则判断 } \mathbf{x} \text{ 为第2类} \end{cases}$