

模式识别

第四章 线性判别函数

郭园方

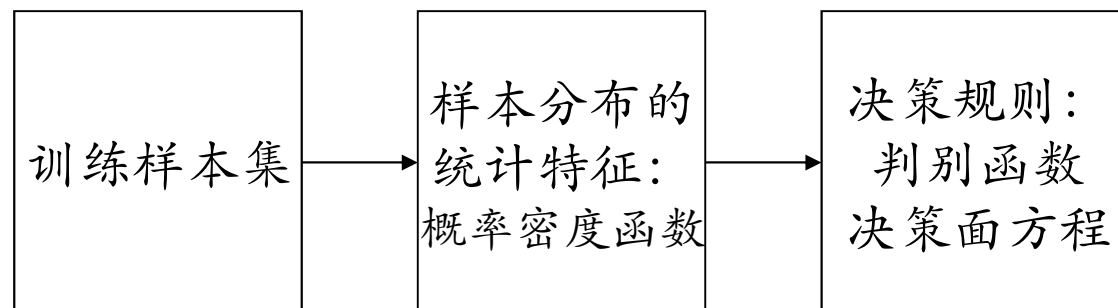
北京航空航天大学计算机学院

内容

- * 引言
- * 线性判别函数的基本概念
- * **Fisher**线性判别函数
- * 感知准则函数
- * 最小平方误差准则函数
- * 多类问题

引言

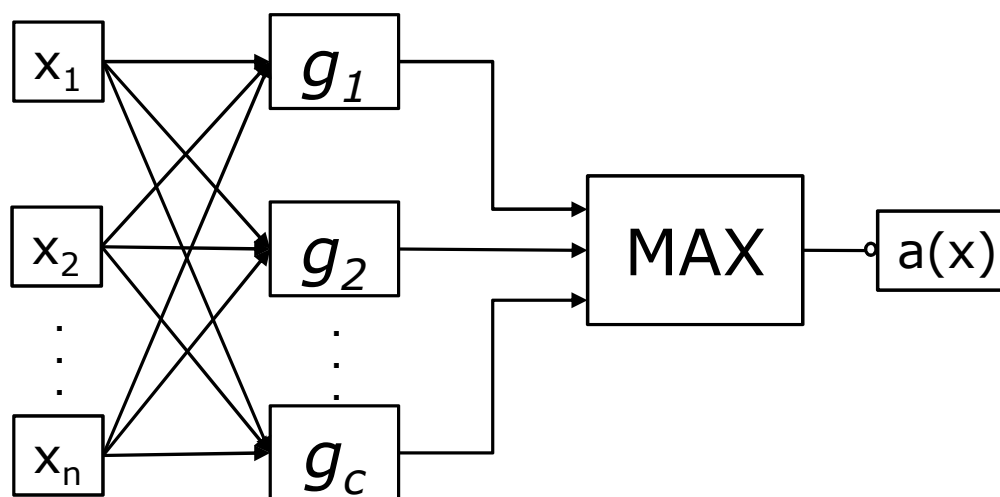
- * 第三章主要讲了类条件概率密度函数的估计
 - * 参数估计方法
 - * 最大似然估计
 - * 贝叶斯估计
 - * 非参数估计方法



引言

- * 基于样本的**Bayes**分类器：通过估计类条件概率密度函数，设计相应的判别函数

分类器
功能结构



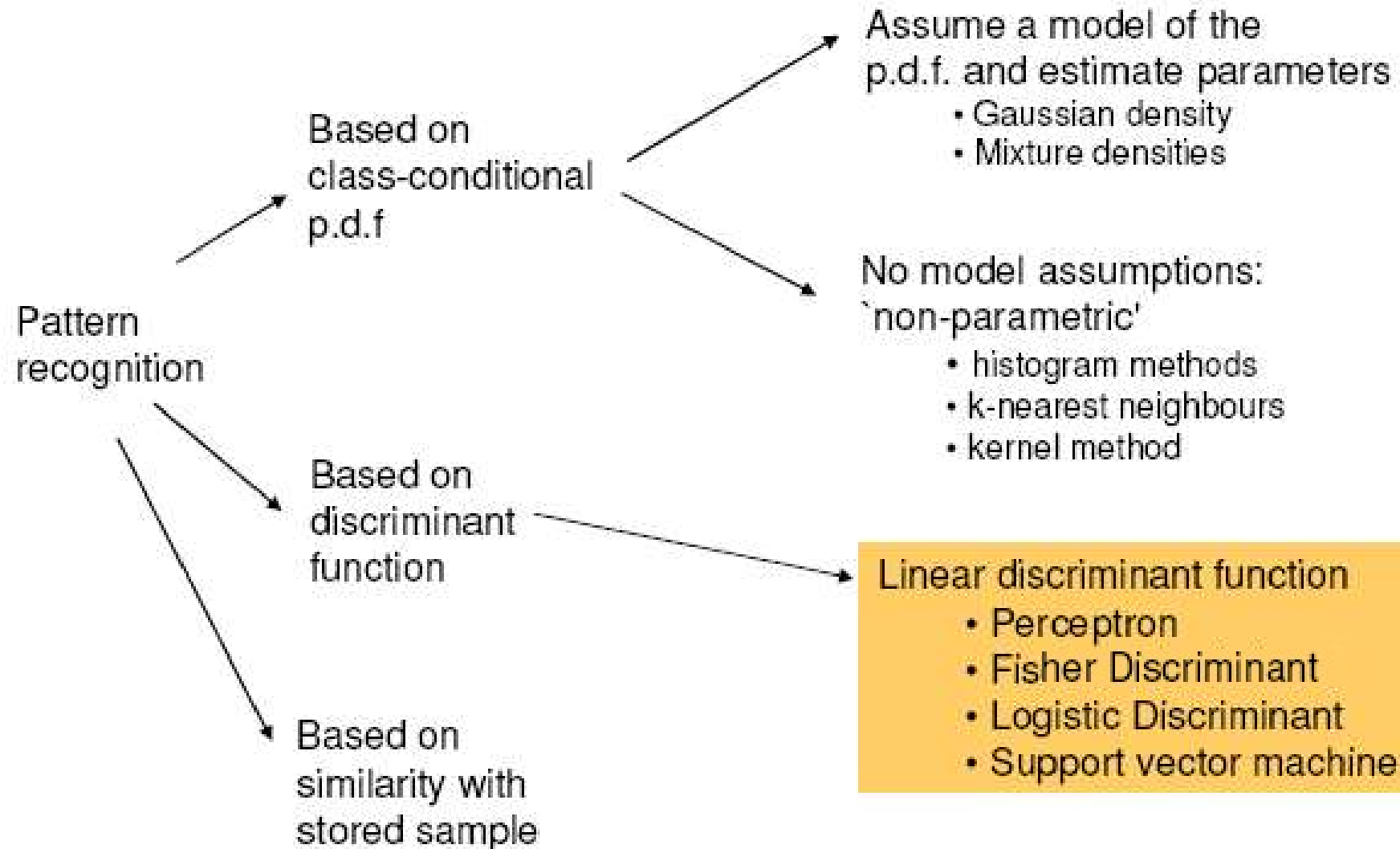
引言

- * 如果可以估计概率密度函数，则可以使用贝叶斯决策来最优的实现分类
- * 最一般情况下适用的“最优”分类器：错误率最小，对分类器设计在理论上具有指导意义。
- * 获取统计分布及其参数很困难，实际问题中并不一定具备获取准确统计分布的条件。
 - * 特征空间维度高
 - * 样本较少
- * 如果估计不了呢？

引言

- * 模式识别还有更多方法!
- * 注意：模式识别（分类）的最终目的是在特征空间中找到分类面!
- * 从样本数据中直接估计参数->根据对样本/问题的理解直接设定判别函数形式，直接求解!
 - * 线性
 - * 非线性

Linear discriminant functions

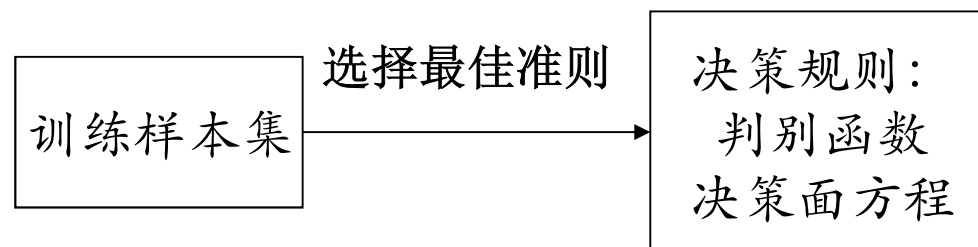


引言

- Parametric Methods
 - Underlying pdfs are known
 - Training samples used to estimate pdf parameters
- Linear Discriminant Functions
 - Forms of discriminant functions are known
 - Similar to non-parametric techniques
 - Sub-optimal, but simple to use

引言

- * 基于样本的直接确定判别函数方法：
 - * 针对各种不同的情况，使用不同的准则函数，设计出满足这些不同准则要求的分类器。
 - * 这些准则的“最优”并不一定与错误率最小相一致：次优分类器。
 - * 实例：正态分布最小错误率贝叶斯分类器在特殊情况下，是线性判别函数 $g(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$ （决策面是超平面），能否基于样本直接确定 \mathbf{w} ？



线性判别函数的基本概念

* d 维空间中的线性判别函数的一般形式:

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

* 其中:

\mathbf{x} 是样本向量-样本在 d 维特征空间中的描述,
 \mathbf{w} 是权向量, w_0 是一个常数(阈值权)。

$$\mathbf{x} = [x_1, x_2, \dots, x_d]^T$$
$$\mathbf{w} = [w_1, w_2, \dots, w_d]^T$$

线性判别函数的基本概念

* 两类问题的分类决策规则:

* 令 $g(\mathbf{x}) = g(x_1) - g(x_2)$

如果 $\begin{cases} g(\mathbf{x}) > 0, & \text{则决策 } \mathbf{x} \in \omega_1 \\ g(\mathbf{x}) < 0, & \text{则决策 } \mathbf{x} \in \omega_2 \\ g(\mathbf{x}) = 0, & \text{可将其任意分类或拒绝} \end{cases}$

线性判别函数的基本概念

$$g(\mathbf{x}) = 0$$

- * 定义了一个决策面，分开了两类
- * 如果 $g(\mathbf{x})$ 是线性的，则决策面是一个超平面
- * 如果 x_1 和 x_2 都在决策面上，则

$$w^T x_1 + w_0 = w^T x_2 + w_0$$

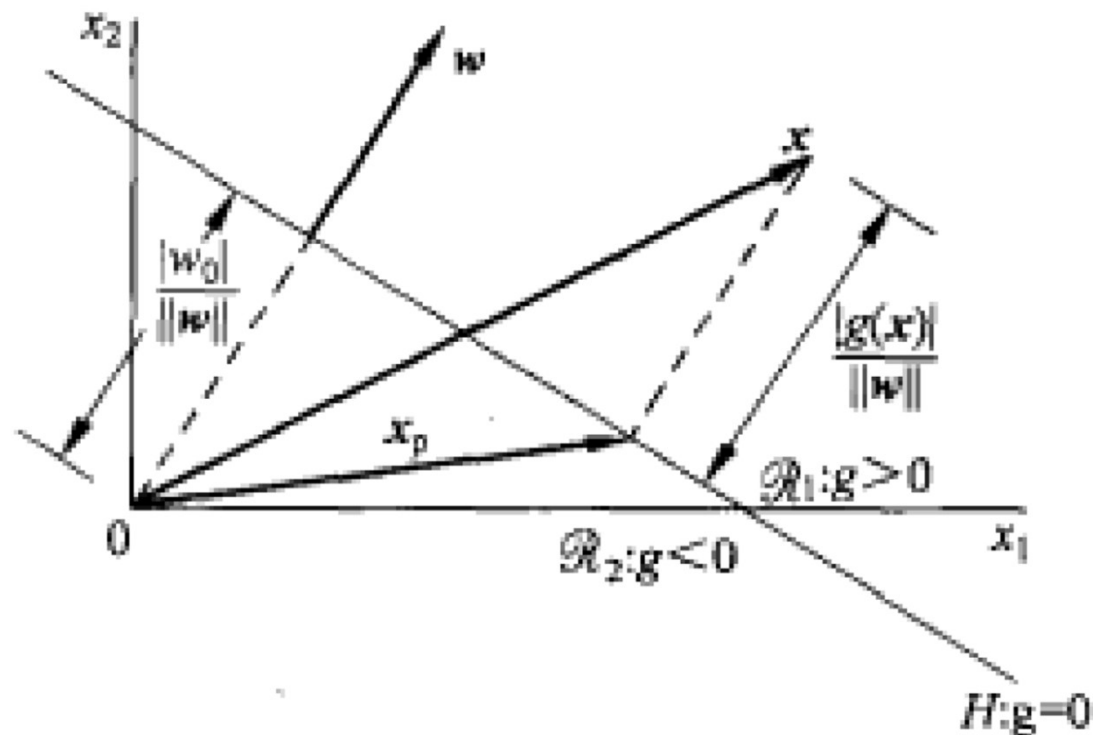
- * 或

$$w^T (x_1 - x_2) = 0$$

线性判别函数的基本概念

* 说明 w 与超平面上任一向量正交 - w 是 H 的法向量

一般来说 H 将空间
分成两个决策域
正侧 vs 负侧



线性判别函数的基本概念

* 判别函数 $g(\mathbf{x})$ 也可以看作是特征空间中某个点 \mathbf{x} 到超平面距离的一种代数度量

* 把 \mathbf{x} 表示为

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

\mathbf{x}_p 是 \mathbf{x} 到 H 的投影向量

r 是 \mathbf{x} 到 H 的垂直距离

$\frac{\mathbf{w}}{\|\mathbf{w}\|}$ 是 \mathbf{w} 方向上的向量

线性判别函数的基本概念

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

* 上式也可写为 $r = \frac{g(\mathbf{x})}{\|\mathbf{w}\|}$

线性判别函数的基本概念

- * 特殊情况 - 若 \mathbf{x} 为原点

$$g(\mathbf{x}) = \omega_0$$

- * 则，原点到超平面的距离为

$$r_0 = \frac{\omega_0}{\|\mathbf{w}\|}$$

- * $\omega_0 > 0$ - 原点在 H 的负侧
- * $\omega_0 < 0$ - 原点在 H 的正侧
- * $\omega_0 = 0$ - 超平面通过原点

线性判别函数的基本概念

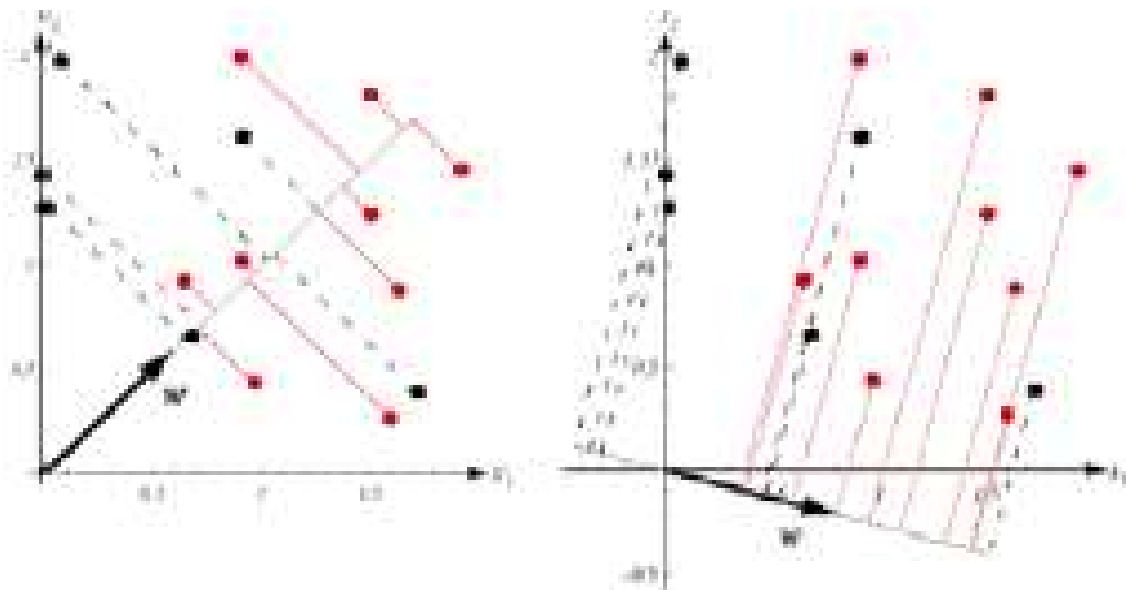
- * 设计线性分类器的主要步骤:
 - * 有一组具有类别标志的样本集
 - * 根据实际情况确定一个准则函数 J ,满足:
 - * J 是样本集和 w, w_0, a 的函数
 - * J 的值能反映分类器的性能, 它的极值解对应于“最好”的决策
 - * 利用最优化方法求出准则函数的极值解和 w, w_0, a , 进而得到 $g(x)$

Fisher线性判别分析

- * 又叫 **Linear discriminant Analysis (LDA)**
- * 1936年的方法
- * **R. A. Fisher**所提出
- * 沿用至今

Fisher线性判别分析

- The Fisher's approach projects d -dimensional data onto a line, w
- It is expected that these projections are well separated by class



Fisher线性判别分析

Fisher判别的基本思想：

希望投影后的一维数据满足：

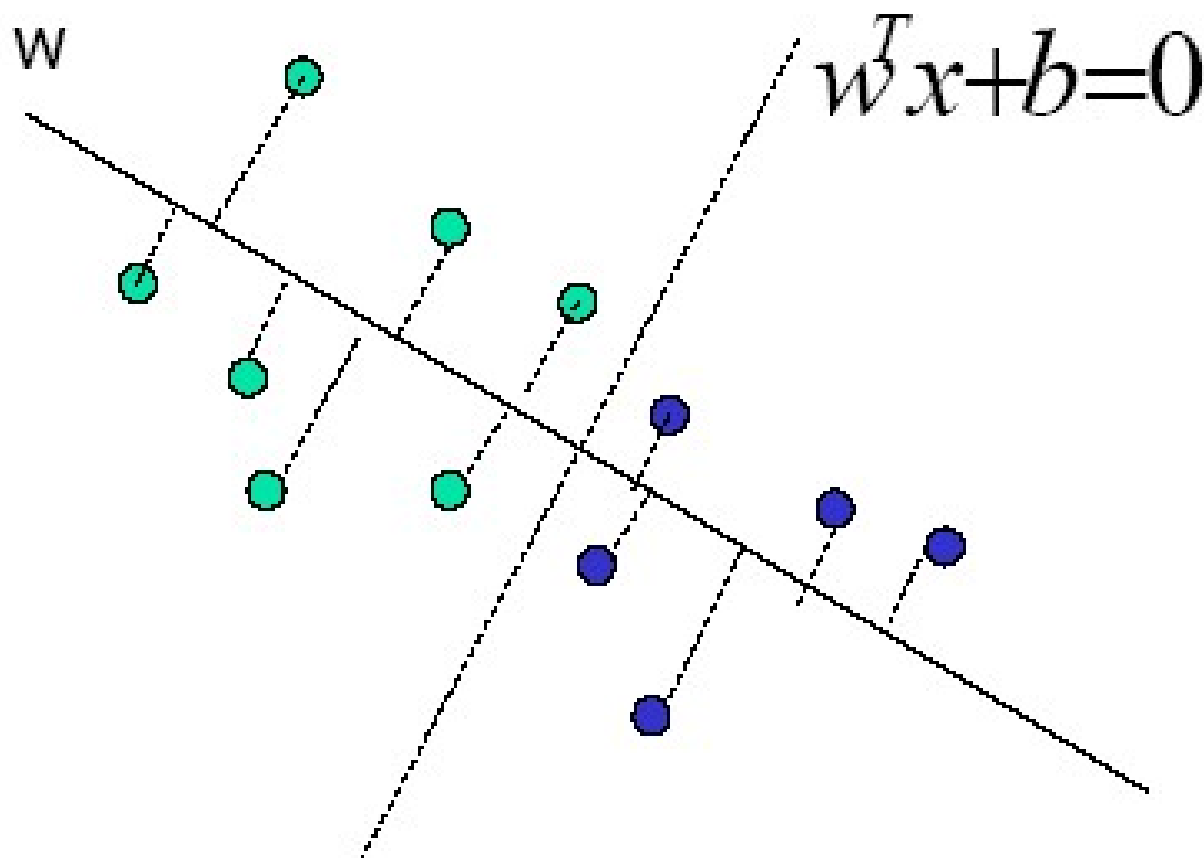
- 两类之间的距离尽可能远；
- 每一类自身尽可能紧凑。

准则的描述：

- 用投影后数据的统计性质——均值和离散度的函数作为判别优劣的标准。

Fisher线性判别分析

- 线性判别函数的几何意义：



Fisher线性判别分析

符号含义:

- m_1, m_2 ——两类（原始）数据的均值向量。
- S_1, S_2 分别表示两类（原始）数据的离散度矩阵。
- μ_1, μ_2 分别表示两类（投影后，一维）数据的均值。
- σ_1, σ_2 分别表示两类（投影后，一维）数据的离散度。

Fisher线性判别分析

- 均值向量和离散度矩阵

$$m_i = \frac{1}{N} \sum x \quad i = 1, 2$$

$$S_i = \sum \left((x - m_i)(x - m_i)^T \right) \quad i = 1, 2$$

Fisher线性判别分析

原始数据与做 w 方向投影后数据统计量之间的关系:

$$\mu_i = w^T m_i,$$

$$i = 1, 2.$$

$$\sigma_i^2 = \sum (w^T x - \mu_i)^2$$

$$= w^T \sum (x - m_i)(x - m_i)^T w$$

$$= w^T S_i w.$$

Fisher线性判别分析

- Fisher准则函数:

$$J_F(w) = \frac{(\mu_1 - \mu_2)^2}{\sigma_1^2 + \sigma_2^2}.$$

类间距
总类内离散度

$$w_{opt} = \arg \max J_F(w).$$

Fisher线性判别分析

- 称 $S_b = (m_1 - m_2)(m_1 - m_2)^T$
类间离散度矩阵。
- 称 $S_t = S_1 + S_2$ 类内总离散度矩阵。

$$J_F(w) = \frac{w^T S_b w}{w^T S_t w}.$$

Fisher线性判别分析

- Fisher准则的合理性:

$J_F(w)$ 只与投影方向有关, 与 w 大小无关——
若 w 是一个最优解, kw 也是最优解, k 是任何不为零的常数。

Fisher线性判别分析

- Fisher最佳投影方向的求解:

- 要求:

$$S_t = S_1 + S_2 \text{ 正定。}$$

否则, 存在投影方向 w , 使得

$$w^T S_t w = 0. \quad \text{所有数据被投影到一点上!}$$

$J_F(w)$ 没有极大值。

Fisher线性判别分析

- 求出最佳投影方向上任何一个 w 即可。

- $J_F(w)$ 有上界，最佳投影方向一定存在！

$$J_F(w) \leq \frac{\lambda(S_b)_{\max}}{\lambda(S_t)_{\min}}.$$

$\lambda(S_w)_{\min}, \lambda(S_b)_{\max}$ 分别是矩阵 S_t, S_b 的最小、最大的特征根。

Fisher线性判别分析

- 一定存在一个最优的 w ，满足：

$$w^T S_t w = 1. \quad \text{因为 } S_t \text{ 正定!}$$

- 无约束最优化： $\max \frac{w^T S_b w}{w^T S_t w}.$

等价于带约束的最优化：

$$\max w^T S_b w$$

$$s.t. \quad w^T S_t w = 1.$$

Fisher线性判别分析

- 带等式约束的最优化，用Lagrange乘子法：

$$L(w, \lambda) = w^T S_b w - \lambda(w^T S_t w - 1).$$

$$\frac{\partial L(w, \lambda)}{\partial w} = S_b w - \lambda S_t w.$$

$$\text{最优解满足: } S_b w_{opt} - \lambda S_t w_{opt} = 0.$$

Fisher线性判别分析

根据类间离散度定义:

$$(m_1 - m_2)(m_1 - m_2)^T w_{opt} = \lambda S_t w_{opt}.$$

注意 $(m_1 - m_2)^T w_{opt}$ 是一个数, 记作 c ,

$$w_{opt} = \frac{c}{\lambda} S_t^{-1} (m_1 - m_2).$$

Fisher线性判别分析

- 只关心投影的方向:

$$\begin{aligned}w_{opt} &= S_t^{-1}(m_1 - m_2) \\ &= (S_1 + S_2)^{-1}(m_1 - m_2).\end{aligned}$$

Fisher线性判别分析

- * 至此，我们还没有解决分类问题，只是将 d 维映射到 1 维，将 d 维分类问题转划为 1 维分类问题，如何分类？
- * 确定阈值

Fisher线性判别分析

- 分类阈值 b 的确定:

- 两类均值的中点:

$$b = \frac{\mu_1 + \mu_2}{2}.$$

- 投影后数据的均值（ n_1, n_2 是两类样本的个数）

$$b = \frac{n_1 \mu_1 + n_2 \mu_2}{n_1 + n_2}.$$

感知准则函数

- * **Perceptron**

- * 感知准则函数是五十年代由**Rosenblatt**提出的一种自学习判别函数生成方法，由于**Rosenblatt**企图将其用于脑模型感知器，因此被称为感知准则函数。其特点是随意确定的判别函数初始值，在对样本分类训练过程中逐步修正直至最终确定。

感知准则函数-基本概念

- * 为了讨论方便，把向量 \mathbf{x} 增加一维，为

$$\mathbf{y} = [1, x_1, x_2, \dots, x_d]^T$$

- * \mathbf{y} 为增广的样本向量

- * 增广的权向量为

$$\mathbf{a} = [w_0, w_1, w_2, \dots, w_d]^T$$

- * 线性判别函数变为

$$g(\mathbf{y}) = \mathbf{a}^T \mathbf{y}$$

感知准则函数-基本概念

- * 线性可分性：训练样本集中的两类样本在特征空间可以用一个线性分界面正确无误地分开。在线性可分条件下，对合适的(广义)权向量 \mathbf{a} 应有：

如果 $\mathbf{y} \in \omega_1$, 则 $\mathbf{a}^T \mathbf{y} > 0$

如果 $\mathbf{y} \in \omega_2$, 则 $\mathbf{a}^T \mathbf{y} < 0$

感知准则函数

- * 线性可分:

- * 模式分类如可用任一个线性函数来划分, 则这些模式就称为线性可分的, 否则就是非线性可分的。
- * 一旦线性函数的系数 w_k 被确定, 这些函数就可用作模式分类的基础。

感知准则函数-基本概念

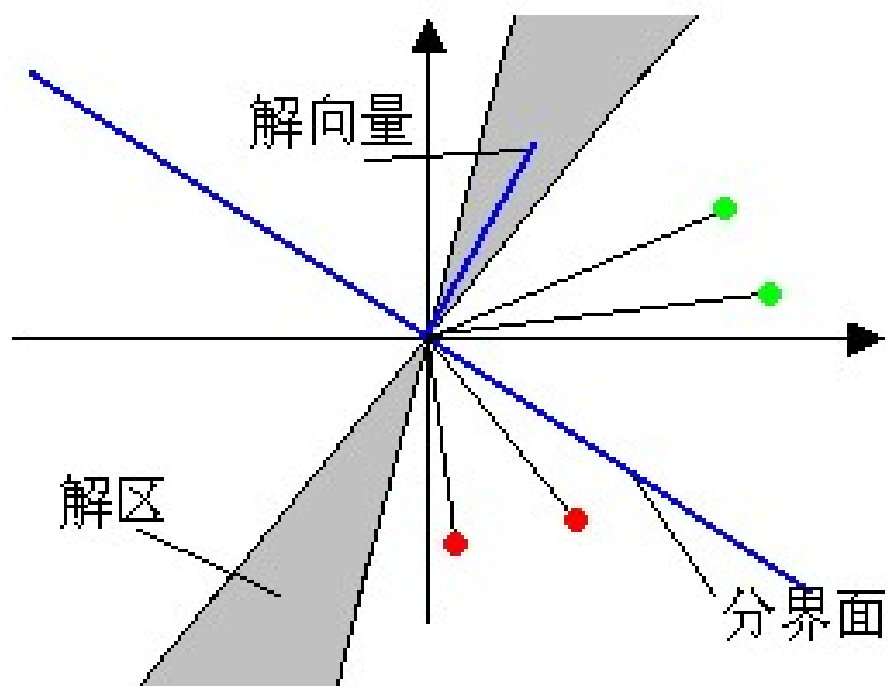
* 规范化样本向量：将第二类样本取其反向向量

$$\mathbf{y}' = \begin{cases} \mathbf{y} & \text{如果 } \mathbf{y} \in \omega_1 \\ -\mathbf{y} & \text{如果 } \mathbf{y} \in \omega_2 \end{cases}$$

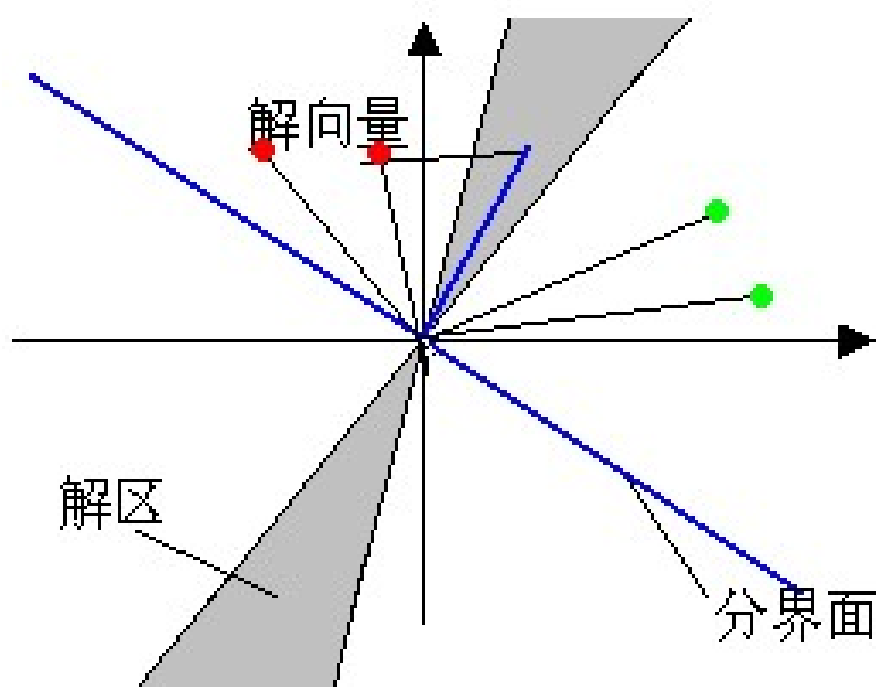


$$\mathbf{a}^T \mathbf{y}'_i > 0 \quad i = 1, \dots, N$$

感知准则函数-解向量与解区



a: 第一类样本
(a) 未规范化



b: 第二类样本
(b) 规范化

感知准则函数

- * 实际中直接用下式会有什么问题么？

$$\mathbf{a}^T \mathbf{y}'_i > 0 \quad i = 1, \dots, N$$

- * 噪声、计算误差可能使得很接近于0的结果出现问题！
- * 怎么办？

感知准则函数

* 引入余量 $b > 0$

$$a^T y'_i > b, i = 1, 2, \dots, N$$

* 如何求解解向量呢?

To find this solution we must first define an objective function $J(a)$

- A good choice is what is known as the **Perceptron** criterion function

$$J_P(a) = \sum_{y \in Y_M} (-a^T y)$$

- where Y_M is the set of examples *misclassified* by a .
- Note that $J_P(a)$ is non-negative since $a^T y < 0$ for all misclassified samples

感知准则函数

To find the minimum of this function, we use gradient descent

- The gradient is defined by

$$\nabla_a J_p(a) = \sum_{y \in Y_M} (-y)$$

- And the gradient descent update rule becomes

$$a(k+1) = a(k) + \eta \sum_{y \in Y_M(k)} y$$

Perceptron rule

- This is known as the **perceptron batch update rule**.
 - The weight vector may also be updated in an “on-line” fashion, this is, after the presentation of each individual example

$$a(k+1) = a(k) + \eta y^{(i)}$$

- where $y^{(i)}$ is an example that has been misclassified by $a(k)$

感知准则函数

- **If classes are linearly separable, the perceptron rule is guaranteed to converge to a valid solution**
 - Some version of the perceptron rule use a variable learning rate $\eta(k)$
 - In this case, convergence is guaranteed only under certain conditions (for details refer to [Duda, Hart and Stork, 2001], pp. 232-235)
- **However, if the two classes are not linearly separable, the perceptron rule will not converge**
 - Since no weight vector a can correctly classify every sample in a non-separable dataset, the corrections in the perceptron rule will never cease
 - One ad-hoc solution to this problem is to enforce convergence by using variable learning rates $\eta(k)$ that approach zero as k approaches infinite

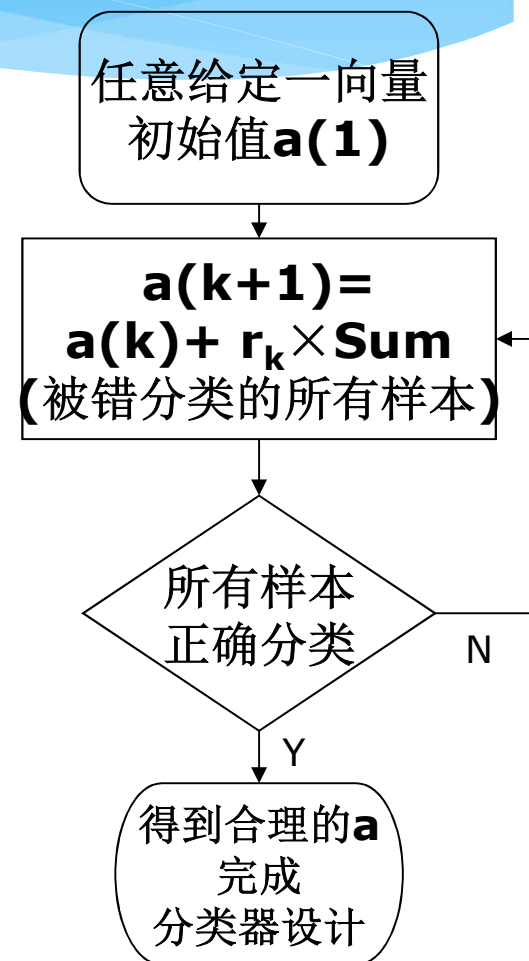
感知准则函数－梯度下降

- * 梯度下降算法：对(迭代)向量沿某函数的负梯度方向修正，可较快到达该函数极小值。

$$\nabla J_p(\mathbf{a}) = \frac{\delta J_p(\mathbf{a})}{\delta \mathbf{a}} = \sum_{\mathbf{y} \in Y^k} (-\mathbf{y})$$
$$\mathbf{a}(k+1) = \mathbf{a}(k) - r_k \nabla J_p(\mathbf{a}) = \mathbf{a}(k) + r_k \sum_{\mathbf{y} \in Y^k} \mathbf{y}$$

感知准则函数-梯度下降

- * 算法具体步骤
- * 初值: 任意给定向量初始值 $\mathbf{a}(1)$
- * 迭代: 第 $k+1$ 次迭代时的权向量 $\mathbf{a}(k+1)$ 等于第 k 次的权向量 $\mathbf{a}(k)$ 加上被错分类的所有样本之和与 r_k 的乘积
- * 终止: 对所有样本正确分类



Fixed-Increment Single-Sample Perceptron Algorithm

- 1 begin initialize \mathbf{a} , $k = 0$
- 2 do $k \leftarrow (k + 1)$
- 3 if $\mathbf{y}^{(k)}$ is misclassified by \mathbf{a}
 then $\mathbf{a} \leftarrow \mathbf{a} + \mathbf{y}^{(k)}$
- 4 until all patterns properly classified
- 5 return
- 6 end

$$\Delta \mathbf{a} = -\eta \nabla J_P(\mathbf{a}) = -\eta \left(\sum_{\mathbf{y} \in Y} (-\mathbf{y}) \right)$$

\mathbf{a} : Modify the weight vector whenever it mis-classifies a single sample

An Example

□ Sample set for two-class case

■ Class 1

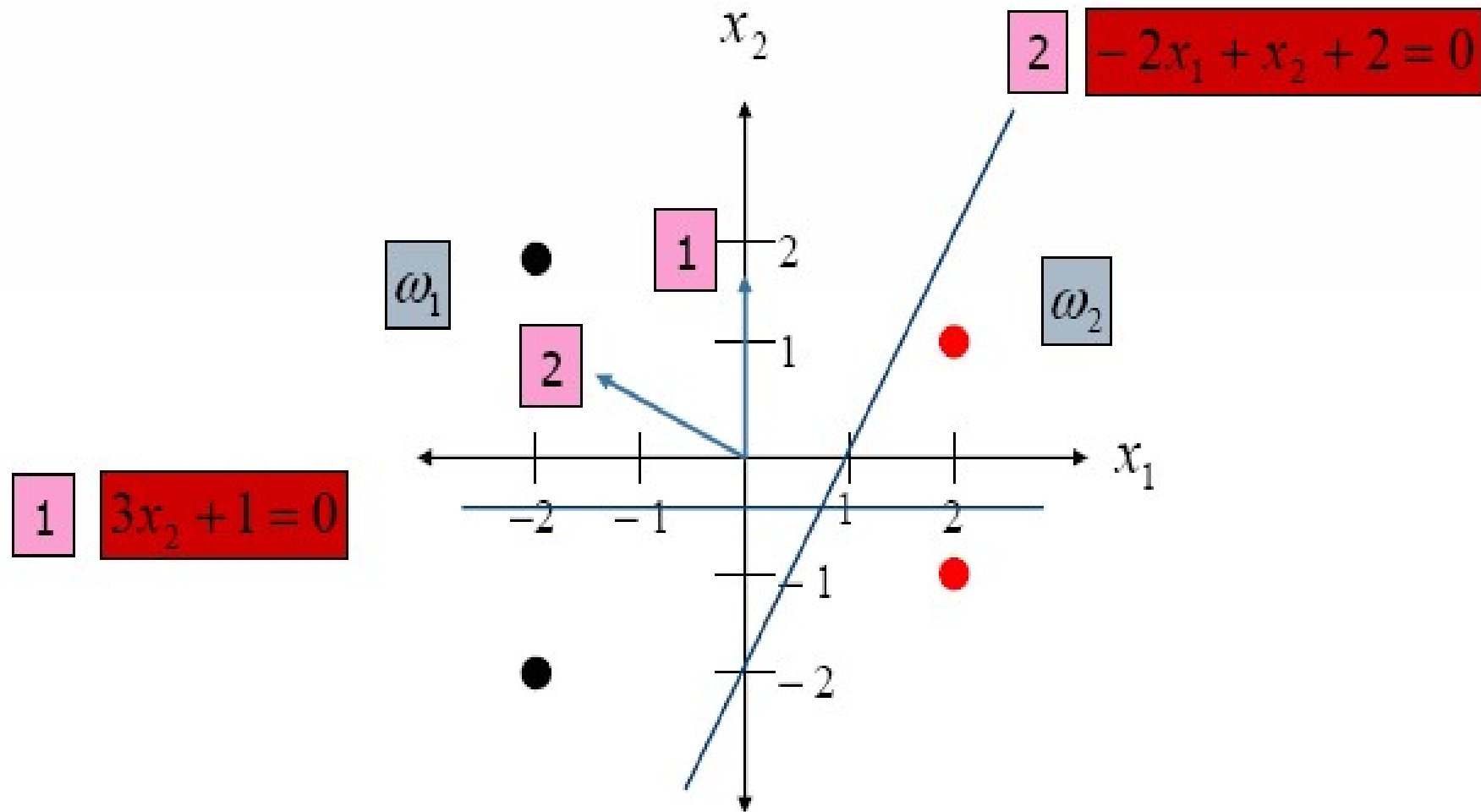
$$\mathbf{x}_1 = \begin{pmatrix} -2 \\ 2 \end{pmatrix}, \mathbf{x}_2 = \begin{pmatrix} -2 \\ -2 \end{pmatrix} \quad \mathbf{y}_1 = \begin{pmatrix} 1 \\ -2 \\ 2 \end{pmatrix}, \mathbf{y}_2 = \begin{pmatrix} 1 \\ -2 \\ -2 \end{pmatrix}$$

■ Class 2

$$\mathbf{x}_3 = \begin{pmatrix} 2 \\ 1 \end{pmatrix}, \mathbf{x}_4 = \begin{pmatrix} 2 \\ -1 \end{pmatrix} \quad \mathbf{y}_3 = \begin{pmatrix} -1 \\ -2 \\ -1 \end{pmatrix}, \mathbf{y}_4 = \begin{pmatrix} -1 \\ -2 \\ 1 \end{pmatrix}$$

■ Initial weight vector

$$\mathbf{a}(1) = \begin{pmatrix} 0 \\ 2 \\ 1 \end{pmatrix} \quad \mathbf{a}(1)^t = (0 \quad 2 \quad 1)$$



Notice : $\mathbf{a}^t \mathbf{y} > 0$

感知准则函数 - 梯度下降示例

□ Example (Cont.)

■ Iterative procedure

$$\square^{(1)} \quad \mathbf{y}^{(1)t} = (1 \quad -2 \quad 2)$$

$$\mathbf{a}(1)^t \mathbf{y}^{(1)} = (0 \quad 2 \quad 1) \begin{pmatrix} 1 \\ -2 \\ 2 \end{pmatrix} = -2 < 0$$

$$\mathbf{a}(2)^t = (0 \quad 2 \quad 1) + (1 \quad -2 \quad 2) = (1 \quad 0 \quad 3)$$

感知准则函数－梯度下降示例

□ Example (Cont.)

■ Iterative procedure

$$\square(2) \quad \mathbf{y}^{(2)t} = (1 \quad -2 \quad -2)$$

$$\mathbf{a}(2)^t \mathbf{y}^{(2)} = (1 \quad 0 \quad 3) \begin{pmatrix} 1 \\ -2 \\ -2 \end{pmatrix} = -5 < 0$$

$$\mathbf{a}(3)^t = (1 \quad 0 \quad 3) + (1 \quad -2 \quad -2) = (2 \quad -2 \quad 1)$$

感知准则函数－梯度下降示例

□ Example (Cont.)

■ Iterative procedure

$$\square(3) \quad \mathbf{y}^{(3)t} = (-1 \quad -2 \quad -1)$$

$$\mathbf{a}(3)^t \mathbf{y}^{(3)} = (2 \quad -2 \quad 1) \begin{pmatrix} -1 \\ -2 \\ -1 \end{pmatrix} = 1 > 0$$

$$\mathbf{a}(4)^t = (2 \quad -2 \quad 1) \text{ (no change)}$$

感知准则函数－梯度下降示例

□ Example (Cont.)

■ Iterative procedure

$$\square(4) \quad \mathbf{y}^{(4)t} = (-1 \quad -2 \quad 1)$$

$$\mathbf{a}(4)^t \mathbf{y}^{(3)} = (2 \quad -2 \quad 1) \begin{pmatrix} -1 \\ -2 \\ 1 \end{pmatrix} = 3 > 0$$

$$\mathbf{a}(5)^t = (2 \quad -2 \quad 1) \quad (\text{no change})$$

感知准则函数－梯度下降示例

□ Example (Cont.)

■ Iterative procedure

$$\square(5) \quad \mathbf{y}^{(5)t} = (1 \quad -2 \quad 2)$$

$$\mathbf{a}(5)^t \mathbf{y}^{(1)} = (2 \quad -2 \quad 1) \begin{pmatrix} 1 \\ -2 \\ 2 \end{pmatrix} = 8 > 0$$

$$\mathbf{a}(6)^t = (2 \quad -2 \quad 1) \quad (\text{no change})$$

感知准则函数－梯度下降示例

□ Example (Cont.)

■ Iterative procedure

$$\square(6) \quad \mathbf{y}^{(6)t} = (1 \quad -2 \quad -2)$$

$$\mathbf{a}(6)^t \mathbf{y}^{(2)} = (2 \quad -2 \quad 1) \begin{pmatrix} 1 \\ -2 \\ -2 \end{pmatrix} = 4 > 0$$

$$\mathbf{a}(7)^t = (2 \quad -2 \quad 1) \text{ (no change)}$$

感知准则函数－梯度下降示例

□ Example (Cont.)

■ Iterative procedure

□(7)

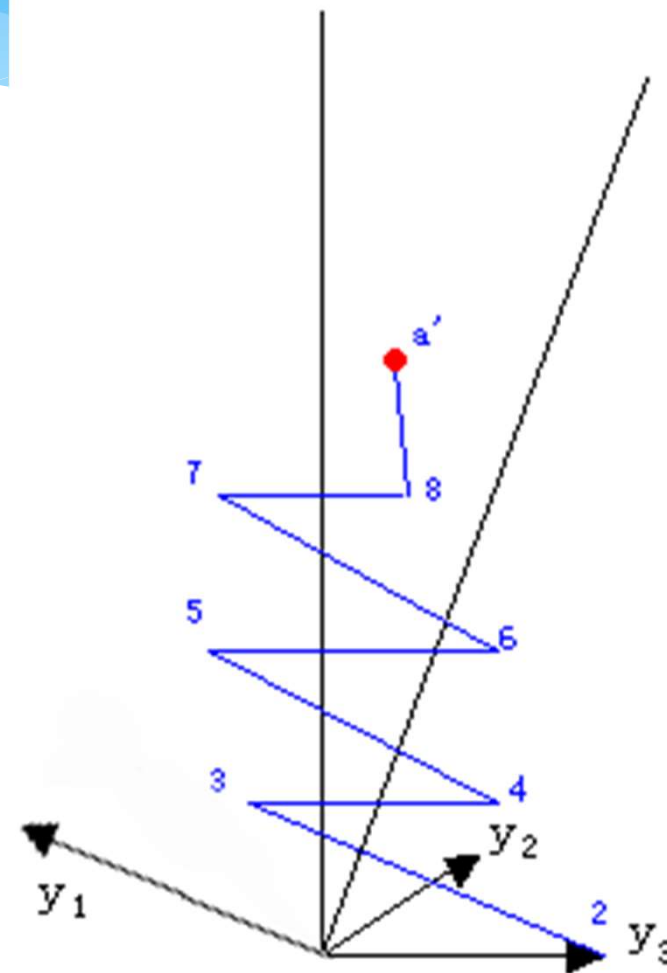
$$\mathbf{y}^{(7)t} = (-1 \quad -2 \quad -1)$$

$$\mathbf{a}^{(7)t} \mathbf{y}^{(7)} = (2 \quad -2 \quad 1) \begin{pmatrix} -1 \\ -2 \\ -1 \end{pmatrix} = 1 > 0$$

$$\mathbf{a}^{(8)t} = (2 \quad -2 \quad 1) \quad (\text{no change})$$

感知准则函数

- * 固定增量法与可变增量法
- * 批量样本修正法与单样本修正法
- * 单样本修正法：样本集视为不断重复出现的序列，逐个样本检查，修正权向量
- * 批量样本修正法：样本成批或全部检查后，修正权向量



感知准则函数-小结

- * 感知准则函数方法的思路:
- * 先随意找一个初始向量 $\mathbf{a}(1)$, 然后用训练样本集中的每个样本来计算。
- * 若发现一个 \mathbf{y} 出现 $\mathbf{a}^T \mathbf{y} < 0$, 则只要 $\mathbf{a}(k+1) = \mathbf{a}(k) + r_k \mathbf{y}$, r_k 为正(步长系数), 则必有 $\mathbf{a}(k+1)^T \mathbf{y} = \mathbf{a}(k)^T \mathbf{y} + r_k \mathbf{y}^T \mathbf{y}$, 就有趋势做到使 $\mathbf{a}(k+1)^T \mathbf{y} > 0$ 。
- * 当然, 修改后的 $\mathbf{a}(k+1)$ 还可以使某些 \mathbf{y} 出现 $\mathbf{a}(k+1)^T \mathbf{y} < 0$ 的情况, 理论证明, 只要训练样本集线性可分, 无论 $\mathbf{a}(1)$ 的初值是什么, 经过有限次叠代, 都可收敛。

最小平方误差准则函数

- * 规范化增广样本向量 \mathbf{y}_i ，增广权向量 \mathbf{a} ，正确分类要求： $\mathbf{a}^T \mathbf{y}_i > 0, i=1, \dots, N$
- * 线性分类器设计 \leftrightarrow 求一组 N 个线性不等式的解

最小平方误差准则函数

- * 样本集增广矩阵 Y 及一组 N 个线性不等式的矩阵表示:

$$Y = \begin{bmatrix} \mathbf{y}_1^T \\ \mathbf{y}_1^T \\ \dots \\ \mathbf{y}_N^T \end{bmatrix} = \begin{bmatrix} y_{11} & y_{12} & \dots & y_{1\hat{d}} \\ y_{21} & y_{22} & \dots & y_{2\hat{d}} \\ \dots & \dots & \dots & \dots \\ y_{N1} & y_{N2} & \dots & y_{N\hat{d}} \end{bmatrix}$$

最小平方误差准则函数

- * 引入余量(目标向量) $\mathbf{b}=[b_1, b_2, \dots, b_N]^T$, b_i 任意给定正常数, $\mathbf{a}^T \mathbf{y}_i = b_i > 0$
- * N 个线性方程的的矩阵表示:

$$Y \mathbf{a} > 0$$

$$Y \mathbf{a} = \mathbf{b}$$

最小平方误差准则函数

- * 定义误差向量 $\mathbf{e} = \mathbf{Y}\mathbf{a} - \mathbf{b}$
- * 定义平方误差准则函数 $J_s(\mathbf{a})$:

$$J_s(\mathbf{a}) = \|\mathbf{e}\|^2 = \|\mathbf{Y}\mathbf{a} - \mathbf{b}\|^2 = \sum_{i=1}^N (\mathbf{a}^T \mathbf{y}_i - b_i)^2$$

- * SSE & MSE

最小平方误差准则函数 函数解

* 最小平方误差准则函数解：

$$\mathbf{a}^* = \underset{\mathbf{a}}{\operatorname{argmin}} J_s(\mathbf{a})$$

MSE方法的思想：对每个样本，设定一个“理想”的判别函数输出值，以最小平方误差为准则求最优权向量

最小平方误差准则函数

MSE准则函数的伪逆解

$$\nabla J_s(\mathbf{a}) = \sum_{i=1}^N 2(\mathbf{a}^T \mathbf{y}_i - b_i) \mathbf{y}_i = 2Y^T (Y\mathbf{a} - \mathbf{b})$$

$$\nabla J_s(\mathbf{a}^*) = 0 \Leftrightarrow Y^T Y \mathbf{a}^* = Y^T \mathbf{b}$$

$$\mathbf{a}^* = (Y^T Y)^{-1} Y^T \mathbf{b} = Y^+ \mathbf{b}$$

Y的
伪逆矩阵

最小平方误差准则函数

MSE方法与其他方法的关系

* 与Fisher方法的关系： 当

$$\mathbf{b} = \begin{bmatrix} N/N_1 \\ \dots \\ N/N_1 \\ N/N_2 \\ \dots \\ N/N_2 \end{bmatrix} \left. \begin{array}{l} \text{ } \\ \text{ } \end{array} \right\} \begin{array}{l} \mathbf{N}_1 \uparrow \\ \mathbf{N}_2 \uparrow \end{array}$$

* MSE解等价于Fisher解

最小平方误差准则函数

MSE方法与其他方法的关系

- * 与Bayes方法的关系：当 $N \rightarrow \infty$ ， $\mathbf{b} = \mathbf{u}_N = [1, 1, \dots, 1]^T$ 时，则它以最小均方误差逼近Bayes判别函数：

$$g(\mathbf{x}) = P(\omega_1 | \mathbf{x}) - P(\omega_2 | \mathbf{x})$$

最小平方误差准则函数

MSE方法的迭代解

- * 伪逆解 $\mathbf{a}^* = \mathbf{Y}^+ \mathbf{b}$, $\mathbf{Y}^+ = (\mathbf{Y}^T \mathbf{Y})^{-1} \mathbf{Y}^T$, 计算量大
- * 实际中常用梯度下降法:

$$\nabla J_s(\mathbf{a}) = \sum_{i=1}^N 2(\mathbf{a}^T \mathbf{y}_i - b_i) \mathbf{y}_i = 2\mathbf{Y}^T (\mathbf{Y}\mathbf{a} - \mathbf{b})$$

$$\begin{cases} \mathbf{a}(1), \text{任意初始化} \\ \mathbf{a}(k+1) = \mathbf{a}(k) - r_k \mathbf{Y}^T (\mathbf{Y}\mathbf{a} - \mathbf{b}) \end{cases}$$

- * 满足一定条件时梯度下降结束

最小平方误差准则函数

MSE方法的迭代解

- * 单样本修正调整权向量
- * Widrow-Hoff算法/最小均方根算法/LMS算法

$$\mathbf{a}(k+1) = \mathbf{a}(k) + r_k(\mathbf{b}_i - \mathbf{a}(k)^T \mathbf{y}_i) \mathbf{y}_i$$

其中 \mathbf{y}_i 是使得 $\mathbf{a}(k)^T \mathbf{y}_i \neq \mathbf{b}_i$ 的样本

Example 1

□ Sample set for two-class case

■ Class 1

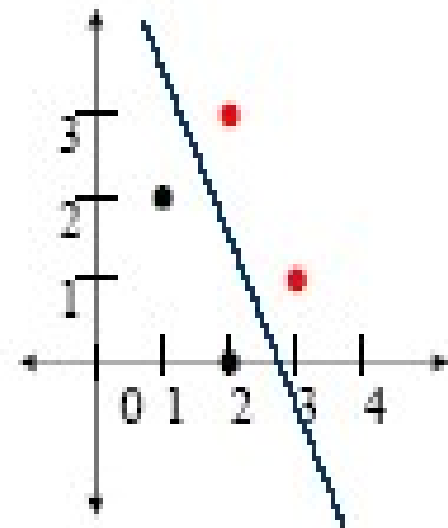
$$\mathbf{x}_1 = \begin{pmatrix} 1 \\ 2 \end{pmatrix}, \quad \mathbf{x}_2 = \begin{pmatrix} 2 \\ 0 \end{pmatrix}$$

Class 2

$$\mathbf{x}_1 = \begin{pmatrix} 3 \\ 1 \end{pmatrix}, \quad \mathbf{x}_2 = \begin{pmatrix} 2 \\ 3 \end{pmatrix}$$

$$\mathbf{Y} = \begin{pmatrix} 1 & 1 & 2 \\ 1 & 2 & 0 \\ -1 & -3 & -1 \\ -1 & -2 & -3 \end{pmatrix}$$

$$\mathbf{Y}^+ = (\mathbf{Y}'\mathbf{Y})^{-1}\mathbf{Y}' = \begin{pmatrix} \frac{5}{3} & \frac{13}{6} & \frac{3}{2} & \frac{7}{6} \\ \frac{4}{3} & \frac{12}{6} & \frac{4}{2} & \frac{12}{6} \\ -\frac{1}{2} & -\frac{1}{6} & -\frac{1}{2} & -\frac{1}{6} \\ 0 & -\frac{1}{3} & 0 & -\frac{1}{3} \end{pmatrix}$$



$$\mathbf{b} = \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \end{pmatrix} \longrightarrow \mathbf{a} = \mathbf{Y}^+\mathbf{b} = \begin{pmatrix} \frac{11}{3} \\ \frac{4}{3} \\ \frac{2}{3} \\ -\frac{1}{3} \end{pmatrix} \longrightarrow \mathbf{a}' \begin{pmatrix} 1 \\ x_1 \\ x_2 \end{pmatrix} = 0 \longrightarrow \frac{11}{3} - \frac{4}{3}x_1 - \frac{2}{3}x_2 = 0 \longrightarrow \frac{4}{11}x_1 + \frac{2}{11}x_2 = 1$$

最小平方误差准则函数

MSE方法的迭代解示例

□ Verification

$$Y a = b \quad \rightarrow$$

$$\begin{pmatrix} 1 & 1 & 2 \\ 1 & 2 & 0 \\ -1 & -3 & -1 \\ -1 & -2 & -3 \end{pmatrix} \begin{pmatrix} 11 \\ -4 \\ -2 \end{pmatrix} = \begin{pmatrix} 3 \\ 3 \\ 3 \\ 3 \end{pmatrix}$$

$$Y = \begin{pmatrix} 1 & 1 & 2 \\ 1 & 2 & 0 \\ -1 & -3 & -1 \\ -1 & -2 & -3 \end{pmatrix} \quad a = \begin{pmatrix} \frac{11}{3} \\ \frac{4}{3} \\ -\frac{2}{3} \end{pmatrix} \quad b = \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \end{pmatrix}$$

最小平方误差准则函数

MSE方法的迭代解示例

□ Sample set for two-class case

■ Class 1

$$\mathbf{x}_1 = \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \mathbf{x}_2 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \mathbf{x}_3 = \begin{pmatrix} 2 \\ 2 \end{pmatrix} \quad \Rightarrow \quad \mathbf{y}_1 = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}, \mathbf{y}_2 = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}, \mathbf{y}_3 = \begin{pmatrix} 1 \\ 2 \\ 2 \end{pmatrix}$$

■ Class 2

$$\mathbf{x}_1 = \begin{pmatrix} -1 \\ 0 \end{pmatrix}, \mathbf{x}_2 = \begin{pmatrix} -1 \\ 1 \end{pmatrix}, \mathbf{x}_3 = \begin{pmatrix} -1 \\ 2 \end{pmatrix} \quad \Rightarrow \quad \mathbf{y}_1 = \begin{pmatrix} -1 \\ 1 \\ 0 \end{pmatrix}, \mathbf{y}_2 = \begin{pmatrix} -1 \\ 1 \\ -1 \end{pmatrix}, \mathbf{y}_3 = \begin{pmatrix} -1 \\ 1 \\ -2 \end{pmatrix}$$

Example 2

$$y_1 = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}, y_2 = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}, y_3 = \begin{pmatrix} 1 \\ 2 \\ 2 \end{pmatrix}$$

$$\square \quad (Y'Y)^{-1} = \left(\sum_{i=1}^n y_i y_i' \right)^{-1}$$

$$y_1 = \begin{pmatrix} -1 \\ 1 \\ 0 \end{pmatrix}, y_2 = \begin{pmatrix} -1 \\ 1 \\ -1 \end{pmatrix}, y_3 = \begin{pmatrix} -1 \\ 1 \\ -2 \end{pmatrix}$$

$$\begin{aligned} \sum_{i=1}^n y_i y_i' &= \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} (1 \ 0 \ 0) + \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} (1 \ 1 \ 1) + \begin{pmatrix} 1 \\ 2 \\ 2 \end{pmatrix} (1 \ 2 \ 2) \\ &\quad + \begin{pmatrix} -1 \\ 1 \\ 0 \end{pmatrix} (-1 \ 1 \ 0) + \begin{pmatrix} -1 \\ 1 \\ -1 \end{pmatrix} (-1 \ 1 \ -1) + \begin{pmatrix} -1 \\ 1 \\ -2 \end{pmatrix} (-1 \ 1 \ -2) = \begin{pmatrix} 6 & 0 & 6 \\ 0 & 8 & 2 \\ 6 & 2 & 10 \end{pmatrix} \end{aligned}$$

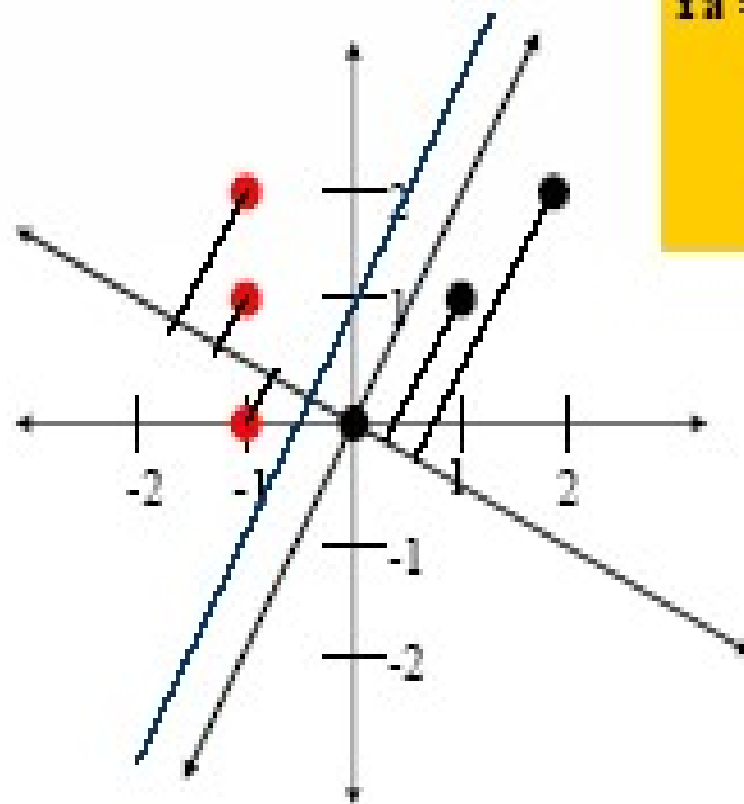
$$\left(\sum_{i=1}^n y_i y_i' \right)^{-1} = \begin{pmatrix} 6 & 0 & 6 \\ 0 & 8 & 2 \\ 6 & 2 & 10 \end{pmatrix}^{-1} \approx \begin{pmatrix} 3 & 0 & 3 \\ 0 & 4 & 1 \\ 3 & 1 & 5 \end{pmatrix}^{-1} \approx \begin{pmatrix} 19 & 3 & -12 \\ 3 & 6 & -3 \\ -12 & -3 & 12 \end{pmatrix}$$

$$\square \quad \mathbf{a} = (\mathbf{Y}'\mathbf{Y})^{-1} \mathbf{Y}'\mathbf{b} = \mathbf{Y}^+ \mathbf{b}$$

$$\mathbf{Y}'\mathbf{b} = \begin{pmatrix} 1 & 1 & 1 & -1 & -1 & -1 \\ 0 & 1 & 2 & 1 & 1 & 1 \\ 0 & 1 & 2 & 0 & -1 & -2 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 0 \\ 6 \\ 0 \end{pmatrix} \quad (\mathbf{Y}'\mathbf{Y})^{-1} = \begin{pmatrix} 19 & 3 & -12 \\ 3 & 6 & -3 \\ -12 & -3 & 12 \end{pmatrix}$$

$$\mathbf{a} = (\mathbf{Y}'\mathbf{Y})^{-1} \mathbf{Y}'\mathbf{b} = \begin{pmatrix} 19 & 3 & -12 \\ 3 & 6 & -3 \\ -12 & -3 & 12 \end{pmatrix} \begin{pmatrix} 0 \\ 6 \\ 0 \end{pmatrix} = \begin{pmatrix} 18 \\ 36 \\ -18 \end{pmatrix} = \begin{pmatrix} 1 \\ 2 \\ -1 \end{pmatrix}$$

$$\mathbf{a}'\mathbf{y} = 0 \quad (1 \ 2 \ -1) \begin{pmatrix} 1 \\ x_1 \\ x_2 \end{pmatrix} = 0 \quad 2x_1 - x_2 + 1 = 0$$



$$2x_1 - x_2 + 1 = 0$$

$$Y\mathbf{a} = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 1 & 1 \\ 1 & 2 & 2 \\ -1 & 1 & 0 \\ -1 & 1 & -1 \\ -1 & 1 & -2 \end{pmatrix} \begin{pmatrix} 1 \\ 2 \\ -1 \end{pmatrix} = \begin{pmatrix} 1 \\ 2 \\ 3 \\ 1 \\ 2 \\ 3 \end{pmatrix}$$

□ Verification

$$\mathbf{Y}\mathbf{a} = \mathbf{b}$$

$$\mathbf{Y}\mathbf{a} = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 1 & 1 \\ 1 & 2 & 2 \\ -1 & 1 & 0 \\ -1 & 1 & -1 \\ -1 & 1 & -2 \end{pmatrix} \begin{pmatrix} 1 \\ 2 \\ -1 \end{pmatrix} = \begin{pmatrix} 1 \\ 2 \\ 3 \\ 1 \\ 2 \\ 3 \end{pmatrix}$$

$$\mathbf{Y} = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 1 & 1 \\ 1 & 2 & 2 \\ -1 & 1 & 0 \\ -1 & 1 & -1 \\ -1 & 1 & -2 \end{pmatrix}$$

$$\mathbf{a} = \begin{pmatrix} 1 \\ 2 \\ -1 \end{pmatrix} \quad \mathbf{b} = \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \end{pmatrix}$$

最小平方误差准则函数

Perceptron Vs. MSE

□ Sample set for two-class case

■ Class 1

$$\mathbf{x}_1 = \begin{pmatrix} -2 \\ 2 \end{pmatrix}, \mathbf{x}_2 = \begin{pmatrix} -2 \\ -2 \end{pmatrix}$$

$$\mathbf{y}_1 = \begin{pmatrix} 1 \\ -2 \\ 2 \end{pmatrix}, \mathbf{y}_2 = \begin{pmatrix} 1 \\ -2 \\ -2 \end{pmatrix}$$

■ Class 2

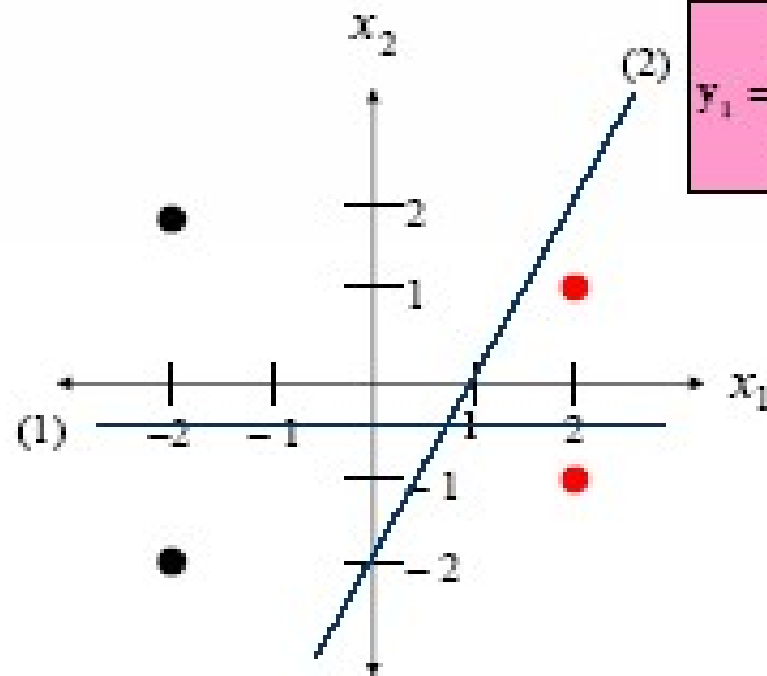
$$\mathbf{x}_3 = \begin{pmatrix} 2 \\ 1 \end{pmatrix}, \mathbf{x}_4 = \begin{pmatrix} 2 \\ -1 \end{pmatrix}$$

$$\mathbf{y}_3 = \begin{pmatrix} -1 \\ -2 \\ -1 \end{pmatrix}, \mathbf{y}_4 = \begin{pmatrix} -1 \\ -2 \\ 1 \end{pmatrix}$$

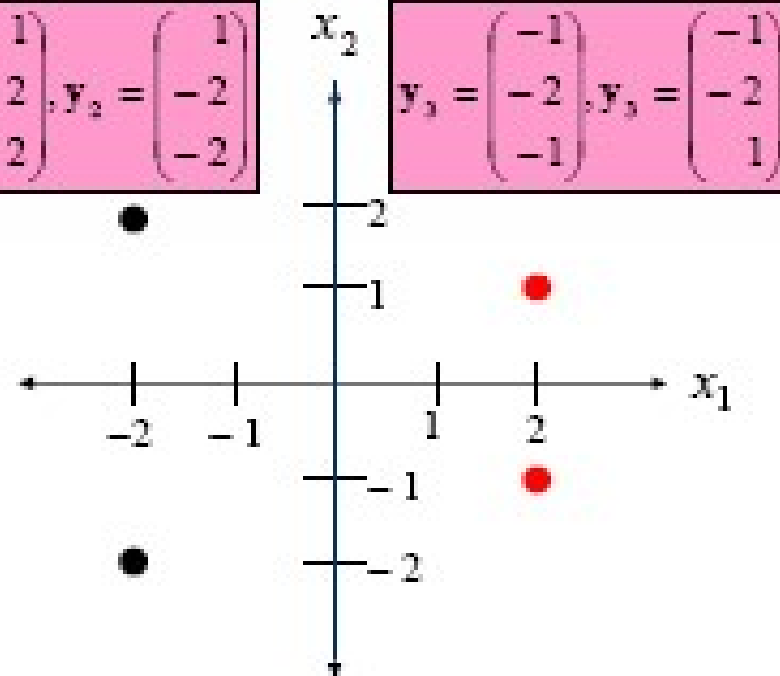
最小平方误差准则函数

Perceptron

MSE



$$y_1 = \begin{pmatrix} 1 \\ -2 \\ 2 \end{pmatrix}, y_2 = \begin{pmatrix} 1 \\ -2 \\ -2 \end{pmatrix}$$



$$y_3 = \begin{pmatrix} -1 \\ -2 \\ -1 \end{pmatrix}, y_4 = \begin{pmatrix} -1 \\ -2 \\ 1 \end{pmatrix}$$

$$Y a = b$$

$$\begin{pmatrix} 1 & -2 & 2 \\ 1 & -2 & -2 \\ -1 & -2 & -1 \\ -1 & -2 & 1 \end{pmatrix} \begin{pmatrix} 2 \\ -2 \\ 1 \end{pmatrix} = \begin{pmatrix} 8 \\ 4 \\ 1 \\ 3 \end{pmatrix}$$

$$a = \begin{pmatrix} 2 \\ -2 \\ 1 \end{pmatrix}$$

$$\begin{pmatrix} 1 & -2 & 2 \\ 1 & -2 & -2 \\ -1 & -2 & -1 \\ -1 & -2 & 1 \end{pmatrix} \begin{pmatrix} 0 \\ 0 \\ -20 \\ 0 \end{pmatrix} = \begin{pmatrix} 40 \\ 40 \\ 40 \\ 40 \end{pmatrix}$$

$$a = \begin{pmatrix} 0 \\ -20 \\ 0 \end{pmatrix}$$

40

最小平方误差准则函数

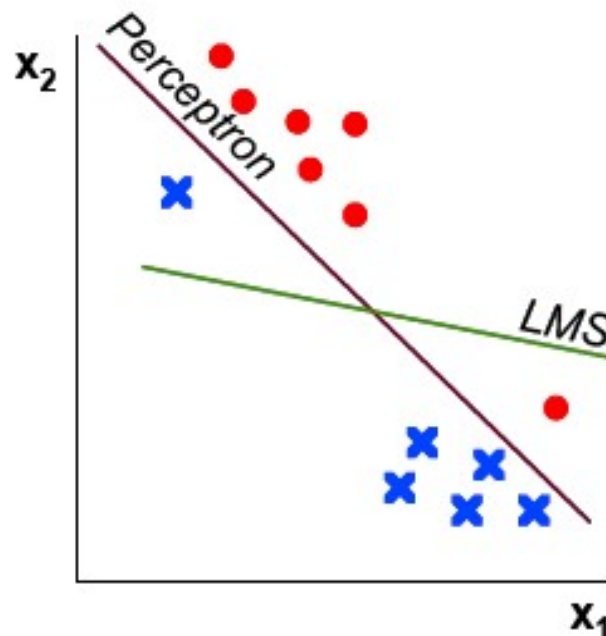
Perceptron Vs MSE

■ Perceptron rule

- The perceptron rule always finds a solution if the classes are linearly separable, but does not converge if the classes are non-separable

■ MSE criterion

- The MSE solution has guaranteed convergence, but it may not find a separating hyperplane if classes are linearly separable
 - Notice that MSE tries to minimize the sum of the squares of the distances of the training data to the separating hyperplane, as opposed to finding this hyperplane



多类问题

□ Three approaches

■ $c-1$ two-class problems

□ Separate points assigned to ω_i from those not assigned to ω_i .

■ $c(c-1)/2$ two-class problems

□ Separate every pair of classes.

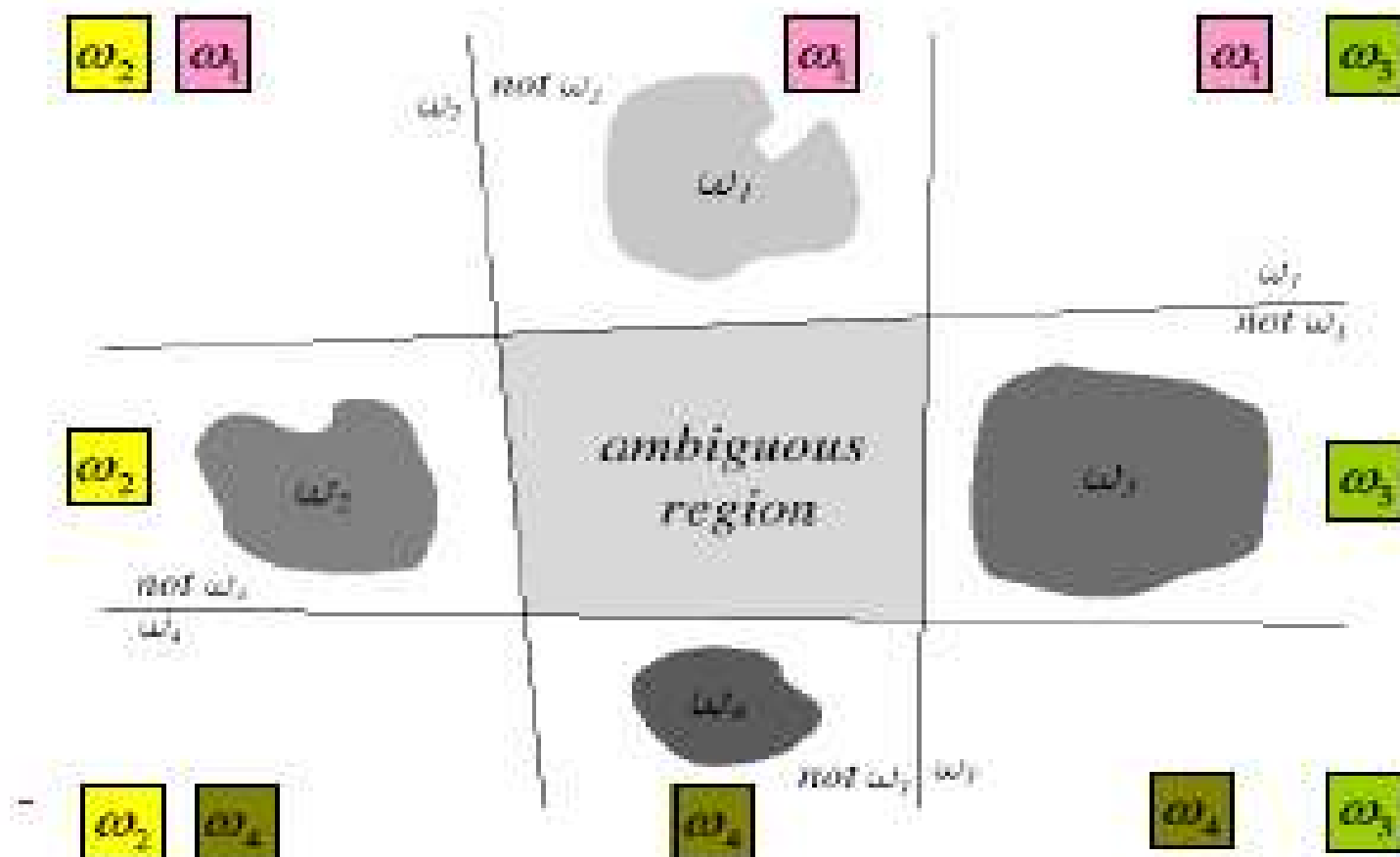
■ Define c linear discriminant functions

□ Assign x to ω_i if $g_i(\mathbf{x}) > g_j(\mathbf{x})$ for all $j \neq i$

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

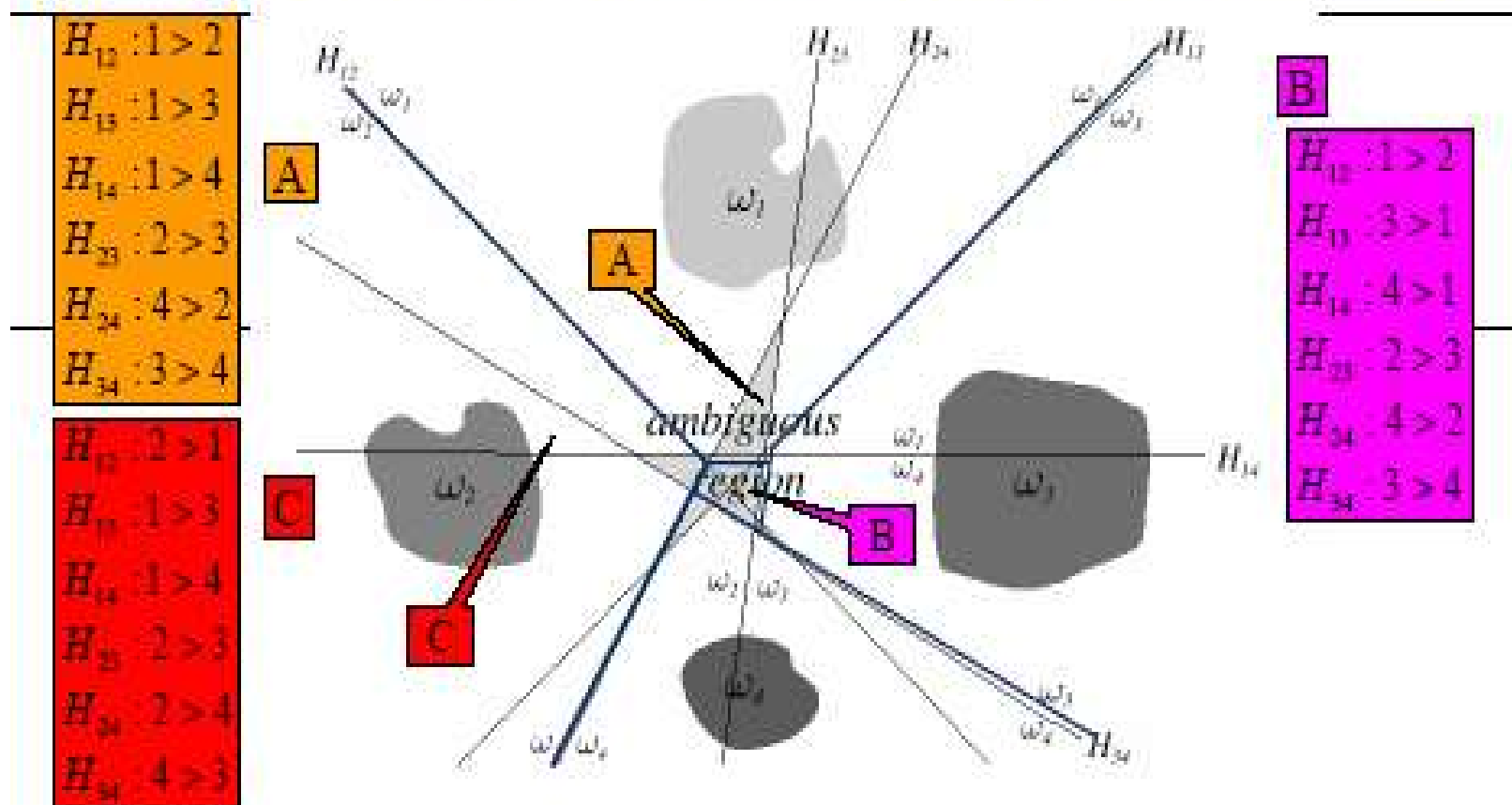
多类问题

Ambiguous region in $c-1$ two-class problems



多类问题

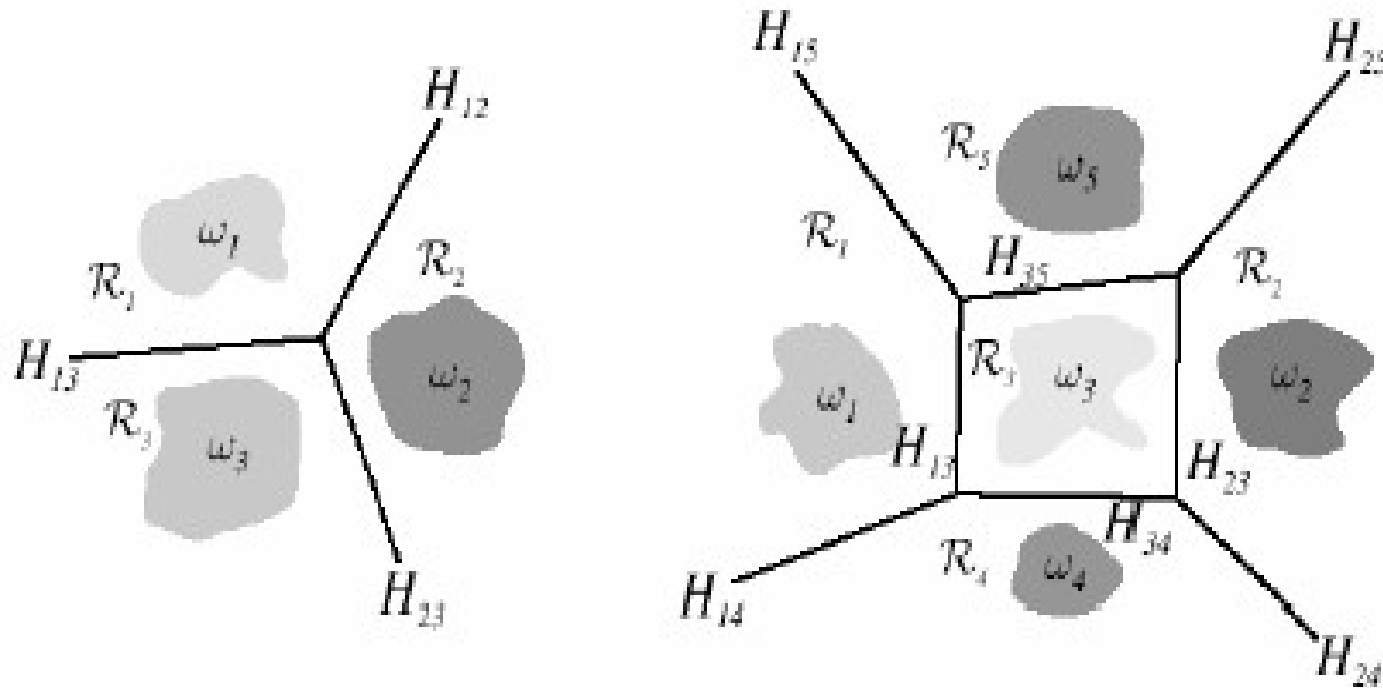
Ambiguous region in $c(c-1)/2$ two-class problems



多类问题

Linear machine

Use c linear discriminant functions



决策树

谢谢