

数据挖掘与机器学习

潘斌

panbin@nankai.edu.cn

范孙楼227

1

上节回顾

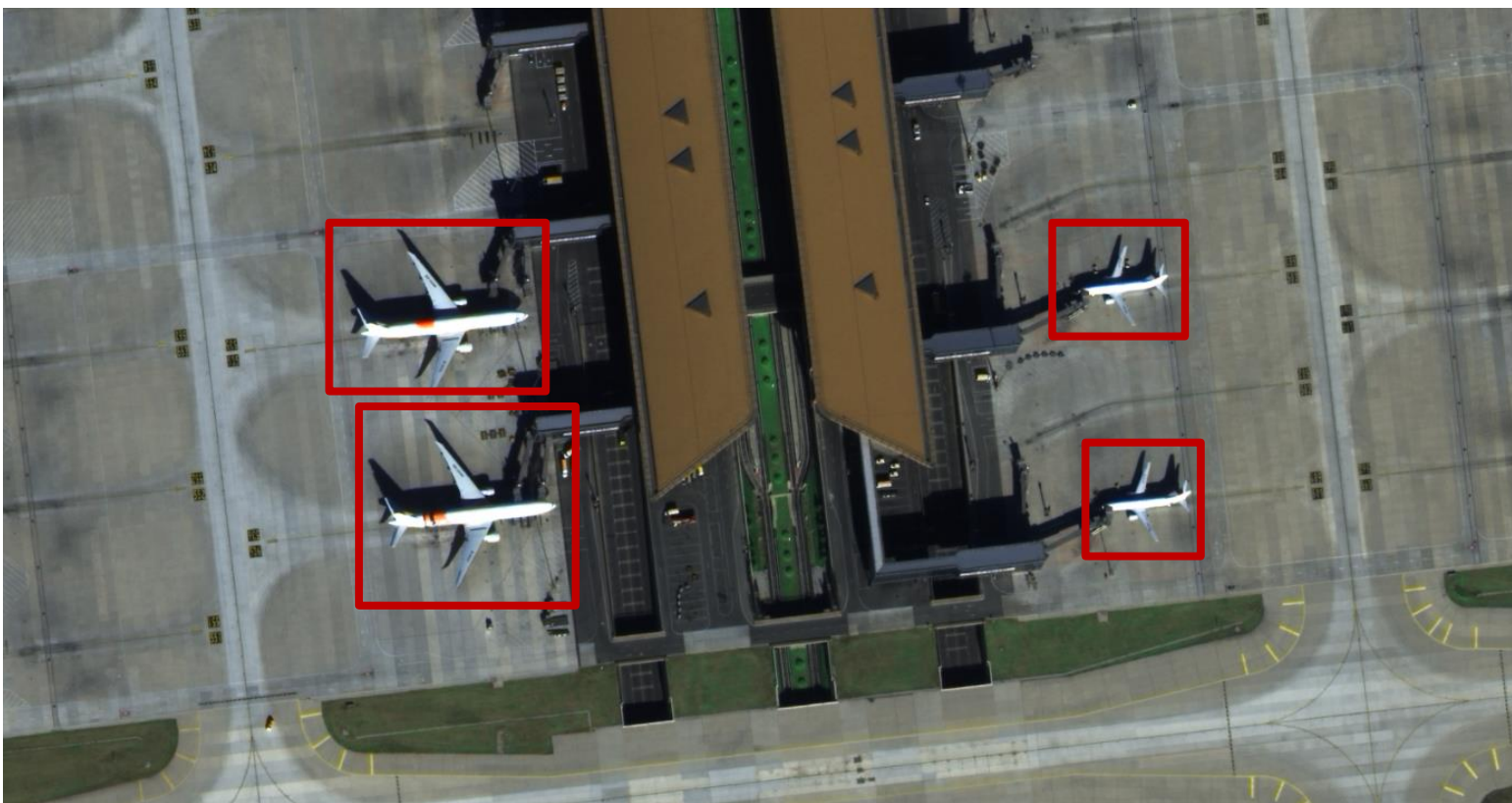
- 线性分类器
 - 垂直平分分类器
 - 感知准则
- 梯度下降

本节提要

- 最小错分样本数准则
- 最小误差准则
- 贝叶斯分类器
 - 最小错误率Bayes决策
 - 最小风险Bayes决策
 - 最大最小Bayes决策
 - 贝叶斯分类器的设计

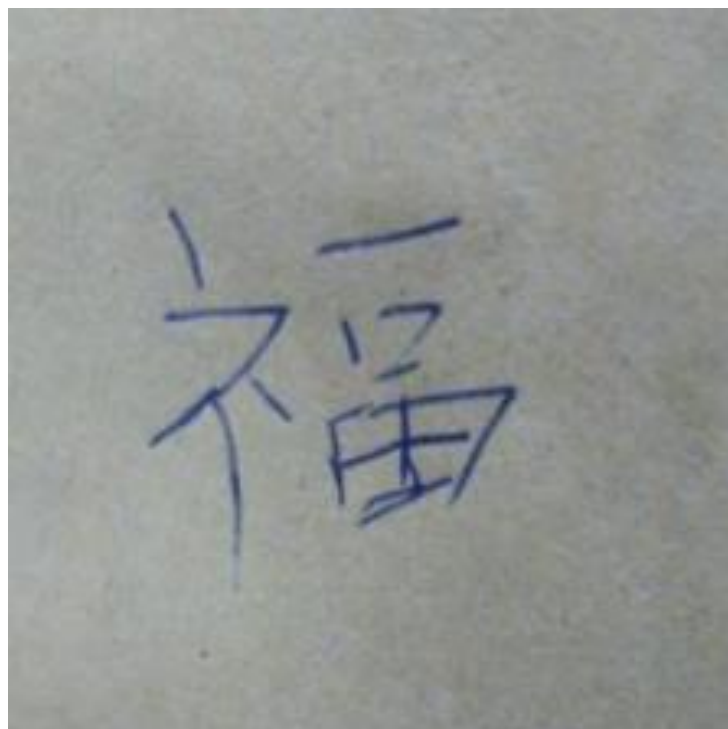
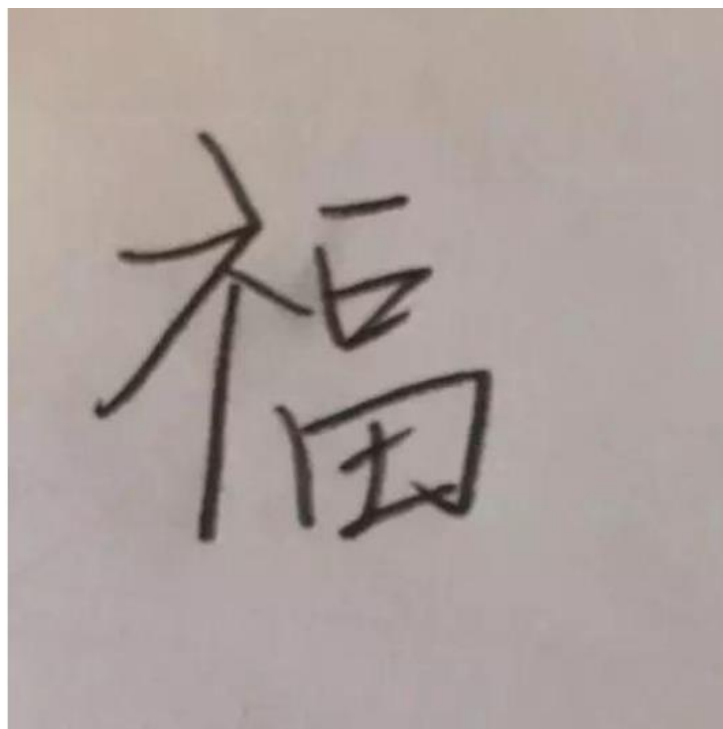
课程大作业（六选一）

■ 任务1：遥感图像飞机检测



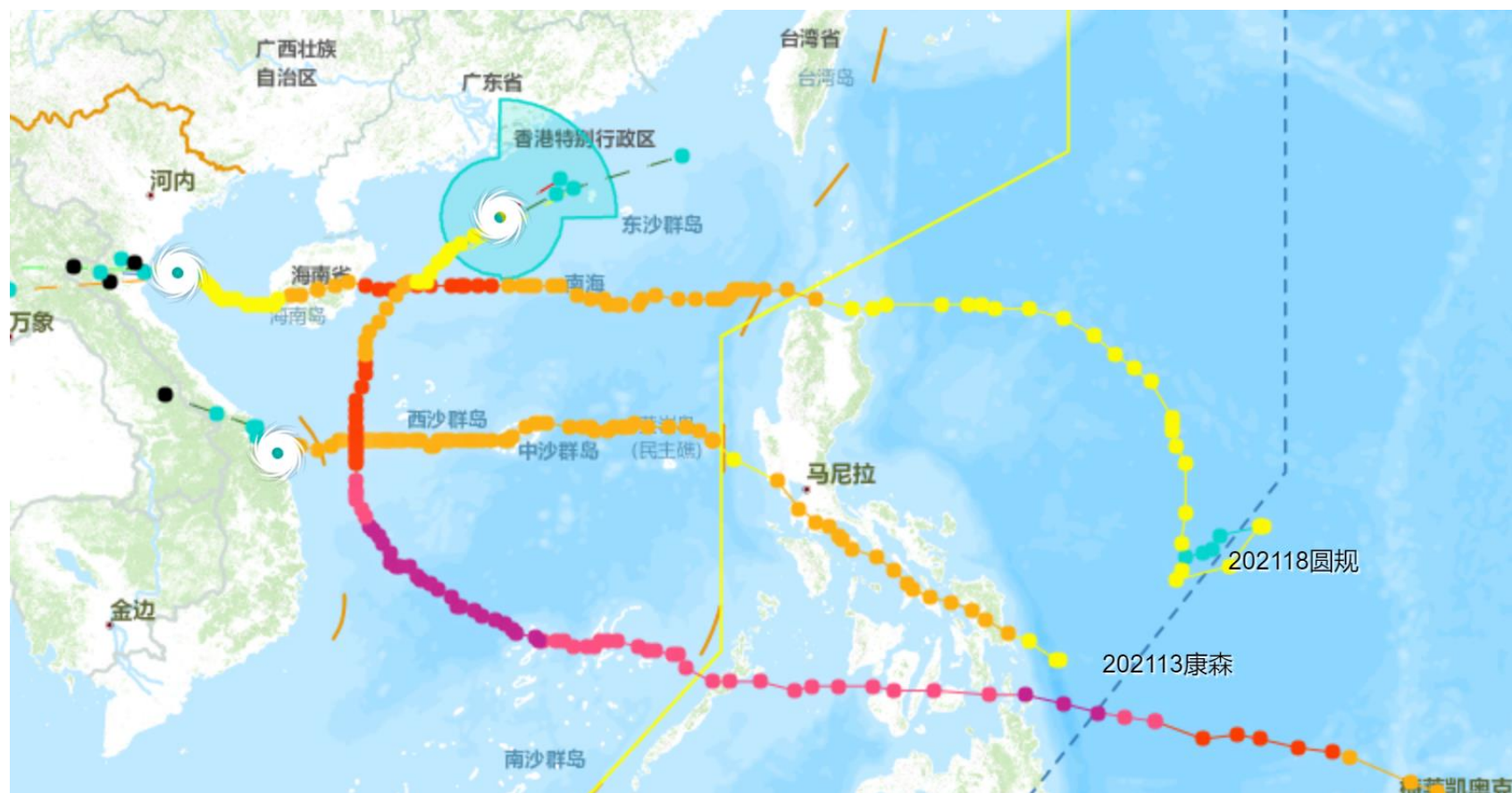
课程大作业（六选一）

■任务2：“福”字识别-解决样本不均衡问题



课程大作业（六选一）

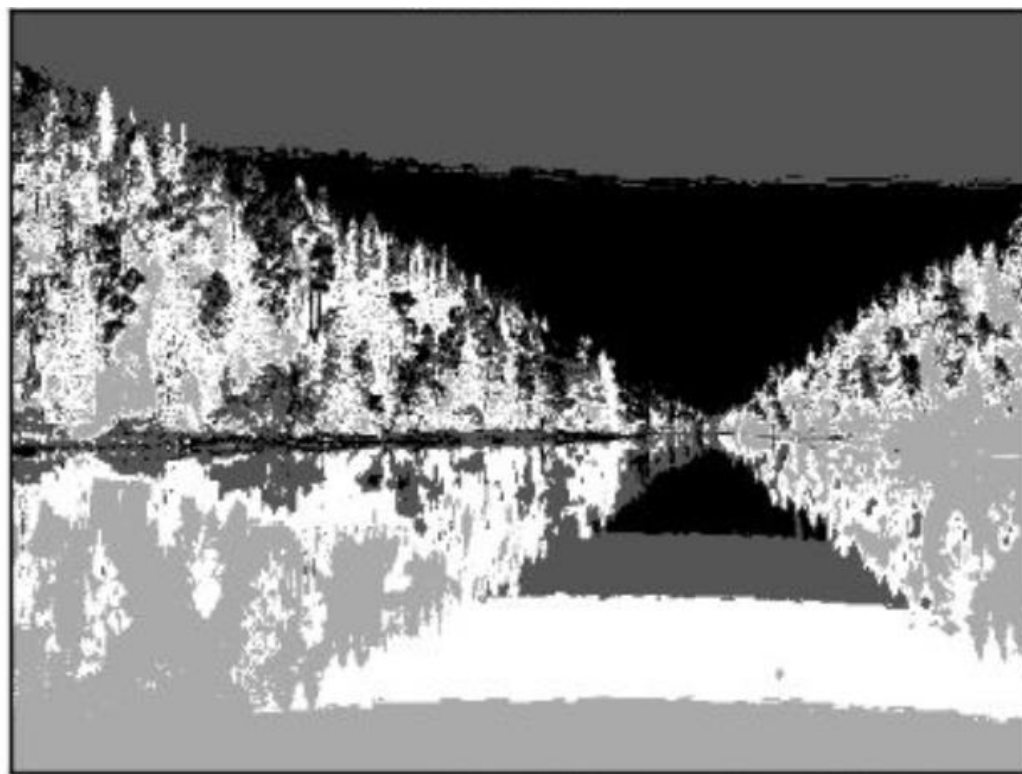
■任务3：台风预报-可以使用RNN/LSTM



66666	1815	21	0018	1815	0	6	LEETPI
2018081012	1	174	1448	1006		13	
2018081018	1	179	1444	1004		15	
2018081100	1	185	1441	1004		15	
2018081106	1	192	1438	1004		15	
2018081112	2	200	1435	1000		18	
2018081118	2	208	1432	995		20	
2018081200	2	217	1427	990		23	
2018081206	2	227	1420	990		23	
2018081212	3	235	1415	982		28	
2018081218	3	245	1407	982		28	
2018081300	3	253	1398	982		28	

课程大作业（六选一）

- 任务4：图像区域分割提取-如何保持空间相关性？



课程大作业（六选一）

■ 任务5：特征选择-特征过多？

Did the blue team get first blood?



电竞比赛结果预测

课程大作业（六选一）

■ 任务6：模型对比- 如何使用不同模型？



手机价格分类

- 任务目标：使用3种课上未提及的模型完成分类任务
- 以下会讲：贝叶斯分类器、线性分类器、决策树、BP神经网络、SVM

5 最小错分样本数准则

- 5.1 问题与思路
- 5.2 最小错分样本数准则一
- 5.3 最小错分样本数准则二
- 5.4 特点

5.1 问题与思路

- 问题的提出
 - 感知准则只适用线性可分样本集——无错分
 - 实际情况未必线性可分——有错分
 - 另外线性可分的判断也很困难
 - 既然存在错分样本——求错分样本数最少

5.1 问题与思路

- 数学描述

- 仿照线性可分样本集的规范化 (ω_2 类样本的增广向量乘以-1)
 - $a^T y_i > 0$ ——正确分类
 - $a^T y_i < 0$ ——错误分类
- 设样本数为N, N个不等式联立
 - $a^T y_i > 0 \quad (i = 1, \dots, N)$
- 求满足不等式最多的解 (权向量)

5.1 问题与思路

- 数学描述
 - 写成矩阵形式

用矩阵形式重写式(4-44)所表示的不等式组,

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

$$Y = \begin{bmatrix} \boldsymbol{y}_1^T \\ \boldsymbol{y}_2^T \\ \vdots \\ \boldsymbol{y}_N^T \end{bmatrix} = \begin{bmatrix} y_{11} & y_{12} & \cdots & y_{1d} \\ y_{21} & y_{22} & \cdots & y_{2d} \\ \cdots & \cdots & \cdots & \cdots \\ y_{N1} & y_{N2} & \cdots & y_{Nd} \end{bmatrix}$$

为使解更可靠,引入余量 $\boldsymbol{b} > 0$

$$Y\boldsymbol{a} \geq \boldsymbol{b} > 0$$

5.2 最小错分样本数准则一

- 最小错分样本数准则一
 - 准则函数

$$\min J_q(\mathbf{a}) = \| (Y\mathbf{a} - \mathbf{b}) - |Y\mathbf{a} - \mathbf{b}| \| ^2$$

- 求极值解
 - 共轭梯度下降法

共轭梯度下降法

在数值线性代数中，共轭梯度法是一种求解对称正定线性方程组 $Ax=b$ 的迭代方法。

事实上，求解 $Ax=b$ 等价于求解： $\min \|Ax - b\|_2^2$ ，将其展开后可以得到： $\min x^T A^T Ax - b^T Ax + b^T b$ ，也就是等价于求解 $\min \frac{1}{2} x^T A^T Ax - b^T Ax$ 。于是解方程问题就转化为了求解二次规划问题(QP)。

共轭梯度法是介于梯度下降法与牛顿法之间的一个方法，是一个**一阶方法**。它克服了梯度下降法收敛慢的缺点，又避免了存储和计算牛顿法所需要的二阶导数信息。

在n维的优化问题中，共轭梯度法最多n次迭代就能找到最优解（是找到，不是接近），但是只针对二次规划问题。

共轭梯度法的思想就是找到n个两两共轭的共轭方向，每次沿着一个方向优化得到该方向上的极小值，后面再沿其它方向求极小值的时候，不会影响前面已经得到的沿哪些方向上的极小值，所以理论上对n个方向都求出极小值就得到了n维问题的极小值。

5.3 最小错分样本数准则二

- 最小错分样本数准则二

- 准则函数

$$\max J_{q2}(\mathbf{a}) = \sum_{i=1}^N \frac{1 + \text{sgn}(y_i \mathbf{a})}{2}$$
$$\text{sgn}(y_i \mathbf{a}) = \begin{cases} +1, & \text{对于 } y_i \mathbf{a} \geq 0 \text{ ①} \\ -1, & \text{对于 } y_i \mathbf{a} < 0 \end{cases}$$

- 求极值解

- 搜索算法

5.4 特点

- 最小错分样本数准则（分类器）的特点
 - 解决两类问题的线性分类器
 - 样本集不限，可以是线性不可分的
 - 求满足不等式个数最多的权向量（最优）
 - 分类器设计过程复杂

6 最小平方误差准则

- 6.1 问题与思路
- 6.2 最小平方误差准则
- 6.3 余量的选择
- 6.4 特点

6.1 问题与思路

- 问题的提出
 - 对于线性不可分问题
 - 最小错分样本数准则——求错分样本数最少
 - 工程上往往是求误差平方和最小

6.1 问题与思路

- 数学描述
 - 引入余量 b_i ，将不等式组改造为等式组
 - $a^T y_i = b_i > 0 \quad (i = 1, \dots, N)$
 - 求满足等式组的最小平方差解（权向量）

6.1 问题与思路

- 数学描述
 - 写成矩阵形式

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

$$Y = \begin{bmatrix} y_1^T \\ y_2^T \\ \vdots \\ y_N^T \end{bmatrix} = \begin{bmatrix} y_{11} & y_{12} & \cdots & y_{1d} \\ y_{21} & y_{22} & \cdots & y_{2d} \\ \cdots & \cdots & \cdots & \cdots \\ y_{N1} & y_{N2} & \cdots & y_{Nd} \end{bmatrix}$$

$$\mathbf{b} = [b_1, b_2, \cdots, b_N]^T$$

6.2 最小平方误差准则

- 最小平方误差准则——工程上常用准则
 - 定义优化准则函数

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

$$J(\mathbf{a}) = \|\mathbf{e}\|^2 = \|Y\mathbf{a} - \mathbf{b}\|^2 = \sum_{n=1}^N (a^T \mathbf{y}_n - b_n)^2$$

6.2 最小平方误差准则

- 最小平方误差准则优化结果
 - 直接求极值解

首先对式(4-63)中的 $J_1(\mathbf{a})$ 求梯度,

$$\nabla J_1(\mathbf{a}) = \sum_{n=1}^N 2(\mathbf{a}^T \mathbf{y}_n - b_n) \mathbf{y}_n = 2Y^T(Y\mathbf{a} - \mathbf{b})$$

令 $\nabla J_1(\mathbf{a}) = 0$, 得

$$Y^T Y \mathbf{a}^* = Y^T \mathbf{b} \quad (4-65)$$

这样, 求解 $Y\mathbf{a} = \mathbf{b}$ 的问题转化为求解 $Y^T Y \mathbf{a}^* = Y^T \mathbf{b}$ 的问题了。这一方程的最大优点是, 矩阵 $Y^T Y$ 是 $d \times d$ 方阵, 而且一般是非奇异的, 因此可唯一地解得

$$\mathbf{a}^* = (Y^T Y)^{-1} Y^T \mathbf{b} = Y^+ \mathbf{b} \quad (4-66)$$

式中 $(d \times N)$ 矩阵

$$Y^+ = (Y^T Y)^{-1} Y^T \quad (4-67)$$

是 Y 的左逆矩阵, \mathbf{a}^* 就是式(4-62)的 MSE 解。

6.3 特点

- 最小平方误差准则（分类器）的特点
 - 解决两类问题的线性分类器
 - 样本集不限，可以是线性不可分的
 - 求最小平方误差的权向量（最优）
 - 分类器设计过程相对简单

Bayes分类器

- 4.1 基本概念
- 4.2 最小错误率Bayes决策
- 4.3 最小风险Bayes决策
- 4.4 最小最大Bayes决策
- 4.5 Bayes分类器设计

4.1 基本概念

- **[错误率]** 几乎所有的分类器在识别时都有可能出现错误分类（简称错分 / 误判）的情况，这种错误分类的可能性称为分类器识别结果的错误概率，简称**错误率 / 误判率**。
- **[正确率]** （通常意义的）正确率 = $1 - \text{错误率}$

4.1 基本概念

- 线性分类器
 - 垂直平分分类器
 - 未经优化，错误率通常较大
 - 感知器
 - 优化（求线性可分样本集的解），最终错误率未知
 - 最小平方误差
 - 优化（样本集MSE的解），最终错误率未知

4.1 基本概念

- **Bayes**分类器设计思路
 - 寻求概率意义上的最小错误率的分类器
 - 即具有最小错分概率的分类器——分类器设计的最优解

4.1 基本概念

- 数学基础回顾
 - 概率论与数理统计
 - 随机事件
 - 概率
 - 条件概率
 - Bayes公式
 - 随机变量
 - 概率密度函数

4.1 基本概念

- 数学基础回顾
 - Bayes分类相关
 - 随机事件——样本的状态/ 类别
 - 概率——状态/ 类别的概率
 - 随机变量——随机向量
 - 概率密度函数

贝叶斯公式

$$P(\theta|X) = \frac{P(X|\theta) \times P(\theta)}{P(X)}$$

$$P(B_i|A) = \frac{P(B_i) P(A|B_i)}{\sum_{j=1}^n P(B_j) P(A|B_j)}$$

- 先验和后验： **$P(\theta)$** 和 **$P(\theta|\mathbf{x})$**

机器学习两大流派——贝叶斯派和频率派

- 频率派旨在求最大似然估计
 - 认为待求参数 θ 是唯一存在的
 - θ 可以是模型参数，也可以是分类标签或预测结果
 - 利用已知的样本结果信息，反推最具有可能（最大概率）导致这些样本结果出现的模型参数值

$$\begin{aligned}\hat{\theta}_{\text{MLE}} &= \arg \max P(X; \theta) \\ &= \arg \max P(x_1; \theta) P(x_2; \theta) \cdots P(x_n; \theta) \\ &= \arg \max \log \prod_{i=1}^n P(x_i; \theta) \\ &= \arg \max \sum_{i=1}^n \log P(x_i; \theta)\end{aligned}$$

$$= \arg \min - \sum_{i=1}^n \log P(x_i; \theta) \quad - \text{负对数似然函数}$$

贝叶斯派和频率派

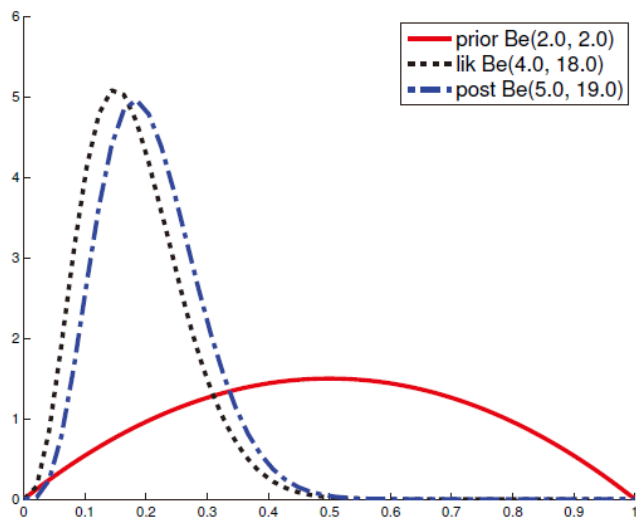
- 贝叶斯派旨在求最大后验估计
 - 认为待求参数 θ 是一个随机变量，符合一定的概率分布
 - 预设一个参数 θ 的概率分布，再用已有样本去修正这个预设（先验概率），得到最有利于样本出现的分布参数（后验概率）

$$\begin{aligned}\hat{\theta}_{\text{MAP}} &= \arg \max P(\theta|X) \\ &= \arg \min -\log P(\theta|X) \\ &= \arg \min -\log P(X|\theta) - \log P(\theta) + \log P(X) \\ &= \arg \min \boxed{-\log P(X|\theta)} - \log P(\theta)\end{aligned}$$

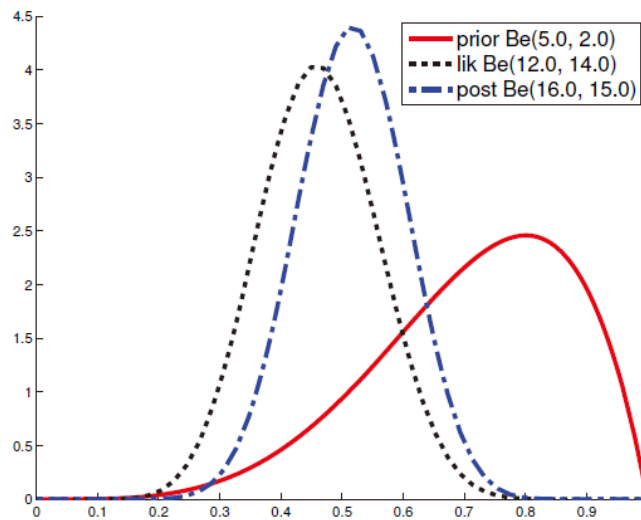
*分号表示确定性，竖条表示不确定性

贝叶斯派和频率派

- 频率派的优势
 - 样本足够大的情况下较容易得到接近无偏的估计
 - 样本少的情况下，偏差较大（例：投5次硬币）
- 贝叶斯派的优势
 - 实际上是基于先验的校正，由于先验的存在，样本少时效果也不会太差
 - 先验非常重要

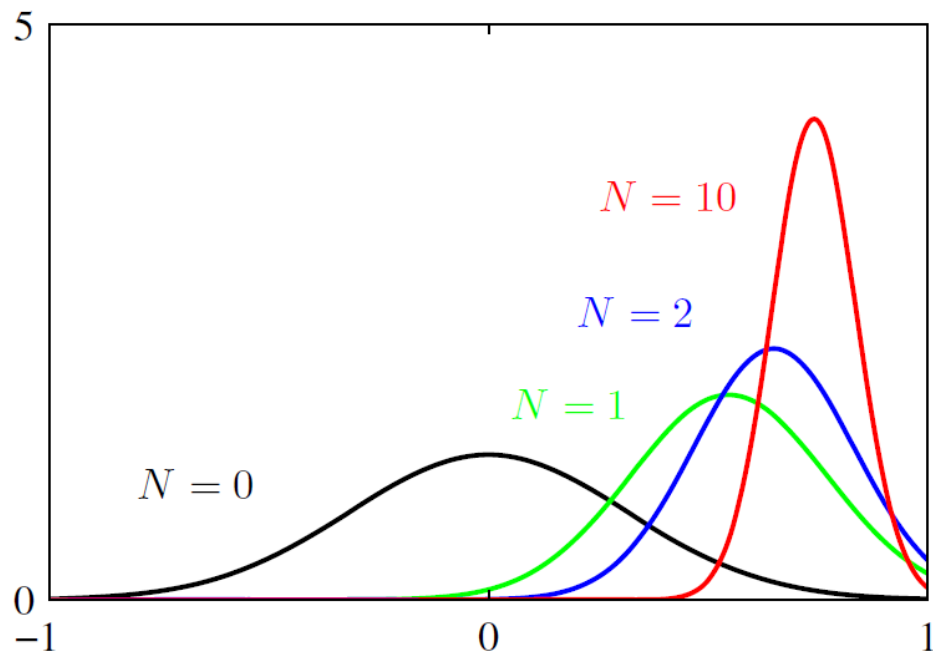


(a)



(b)

- 不同先验的影响



- 训练样本的影响—稀释先验

- 频率派和贝叶斯派的等价关系

$$\mu_N = \frac{\sigma^2}{N\sigma_0^2 + \sigma^2}\mu_0 + \frac{N\sigma_0^2}{N\sigma_0^2 + \sigma^2}\mu_{\text{ML}}$$

$$\frac{1}{\sigma_N^2} = \frac{1}{\sigma_0^2} + \frac{N}{\sigma^2}$$

4.2 最小错误率Bayes决策

- 4.2.1 问题
- 4.2.2 后验概率
- 4.2.3 决策规则
- 4.2.4 实例
- 4.2.5 其它问题
- 4.2.6 特点

4.2.1 问题 (P代表概率, p代表概率密度函数)

- **Bayes**决策已知条件和问题 (先考虑**C = 2**分类问题)
 - 两类问题: ω_1 和 ω_2
 - 先验概率: $P(\omega_1)$ 和 $P(\omega_2)$
 - 类条件概率密度函数: $p(x|\omega_1)$ 和 $p(x|\omega_2)$
 - 发生了一个随机事件, 其观察值为: 特征向量 x
 - 求最小错误率分类器

4. 2. 2 后验概率

- **Bayes**（条件概率）公式
 - $p(x|\omega_1) P(\omega_1) = p(x) P(\omega_1|x)$
 - $p(x|\omega_2) P(\omega_2) = p(x) P(\omega_2|x)$
- 后验概率
 - $P(\omega_1|x) = p(x|\omega_1) P(\omega_1) / p(x)$
 - $P(\omega_2|x) = p(x|\omega_2) P(\omega_2) / p(x)$
 - $p(X) = p(x|\omega_1) P(\omega_1) + p(x|\omega_2) P(\omega_2)$

4. 2. 3 决策规则

- 决策规则

- 比较后验概率，取最大值进行类别判断
- 对于未知样本 \mathbf{x} ，若 $P(\omega_1|\mathbf{x}) > P(\omega_2|\mathbf{x})$ ，则 $\mathbf{x} \in \omega_1$
- 若 $P(\omega_1|\mathbf{x}) < P(\omega_2|\mathbf{x})$ ，则 $\mathbf{x} \in \omega_2$

4.2.3 决策规则

- 等价规则一 后验概率分子
 - 比较分子 $p(x|\omega_1) P(\omega_1)$ 和 $p(x|\omega_2) P(\omega_2)$ ，取最大
 - 对于未知样本 x ，若 $p(x|\omega_1) P(\omega_1) > p(x|\omega_2) P(\omega_2)$ ，则 $x \in \omega_1$
 - 若 $p(x|\omega_1) P(\omega_1) < p(x|\omega_2) P(\omega_2)$ ，则 $x \in \omega_2$

4. 2. 3 决策规则

- 等价规则二 似然比

- 定义似然比函数 $l(x) = p(x|\omega_1) / p(x|\omega_2)$
- 对于未知样本 x , 若 $l(x) > P(\omega_2) / P(\omega_1)$, 则 $x \in \omega_1$
- 若 $l(x) < P(\omega_2) / P(\omega_1)$, 则 $x \in \omega_2$

4. 2. 3 决策规则

- 等价规则三 负对数似然比
 - 定义负对数似然比函数 $h(x) = -\ln l(x) = -\ln p(x|\omega_1) + \ln p(x|\omega_2)$
 - 对于未知样本 x , 若 $h(x) < \ln [P(\omega_1) / P(\omega_2)]$, 则 $x \in \omega_1$
 - 若 $h(x) > \ln [P(\omega_1) / P(\omega_2)]$, 则 $x \in \omega_2$

4. 2. 4 实例

- 已知

- 癌细胞图像识别：正常和异常两类（即 $C = 2$ ）
- 已知未知样本特征观察值： $x = 0.5$
- 又已知 $P(\omega_1) = 0.9$ 和 $P(\omega_2) = 0.1$
- 查函数曲线得 $p(0.5|\omega_1) = 0.2$ 和 $p(0.5|\omega_2) = 0.4$
（可以是先验或学习得到）
- 试对未知样本 $x = 0.5$ 进行分类

4.2.4 实例

- 已知

- 癌细胞图像识别：正常和异常两类（即 $C = 2$ ）
- 已知未知样本特征观察值： $x = 0.5$
- 又已知 $P(\omega_1) = 0.9$ 和 $P(\omega_2) = 0.1$
- 查函数曲线得 $p(0.5|\omega_1) = 0.2$ 和 $p(0.5|\omega_2) = 0.4$
- 试对未知样本 $x = 0.5$ 进行分类

解：利用贝叶斯公式，分别计算出 ω_1 及 ω_2 的后验概率。

$$P(\omega_1|x) = \frac{p(x|\omega_1)P(\omega_1)}{\sum_{j=1}^2 p(x|\omega_j)P(\omega_j)} = \frac{0.2 \times 0.9}{0.2 \times 0.9 + 0.4 \times 0.1} = 0.818$$
$$P(\omega_2|x) = 1 - p(\omega_1|x) = 0.182$$

根据贝叶斯决策规则式(2-2)，有

$$P(\omega_1|x) = 0.818 > P(\omega_2|x) = 0.182$$

所以合理的决策是把 x 归类于正常状态。

4.2.5 最小错误率的说明

- 最小错误率的说明（设 $C=2$, $D=1$ ）
 - 错误率 $P(e)$ 的定义

首先应指出所谓错误率是指平均错误率,以 $P(e)$ 来表示,其定义为

$$P(e) = \int_{-\infty}^{\infty} P(e, \mathbf{x}) d\mathbf{x} = \int_{-\infty}^{\infty} P(e|\mathbf{x}) p(\mathbf{x}) d\mathbf{x} \quad (2-6)$$

其中 $\int_{-\infty}^{\infty} () d\mathbf{x}$ 表示在整个 d 维特征空间上的积分。

对两类别问题,从式(2-2)的决策规则可知,如果 $P(\omega_2|\mathbf{x}) > P(\omega_1|\mathbf{x})$,则决策应为 ω_2 ,显然在作出决策 ω_2 时, \mathbf{x} 的条件错误概率为 $P(\omega_1|\mathbf{x})$;反之,则应为 $P(\omega_2|\mathbf{x})$ 。可表示为

$$P(e|\mathbf{x}) = \begin{cases} P(\omega_1|\mathbf{x}), & \text{当 } P(\omega_2|\mathbf{x}) > P(\omega_1|\mathbf{x}) \\ P(\omega_2|\mathbf{x}), & \text{当 } P(\omega_1|\mathbf{x}) > P(\omega_2|\mathbf{x}) \end{cases} \quad (2-7)$$

4.2.5 最小错误率的说明

- 最小错误率的说明（设**C=2**，**D=1**）

- 错误率**P(e)**的推导

- 做 ω_1 判别时的错误概率

$$P(e_{12}) = \int_{\mathfrak{R}_1} P(\omega_2 | x) p(x) dx = \int_{\mathfrak{R}_1} P(\omega_2) p(x | \omega_2) dx$$

- 做 ω_2 判别时的错误概率

$$P(e_{21}) = \int_{\mathfrak{R}_2} P(\omega_1 | x) p(x) dx = \int_{\mathfrak{R}_2} P(\omega_1) p(x | \omega_1) dx$$

- 总错误概率

$$P(e) = \int_{\mathfrak{R}_1} P(\omega_2) p(x | \omega_2) dx + \int_{\mathfrak{R}_2} P(\omega_1) p(x | \omega_1) dx$$

4.2.5 最小错误率的说明

- 最小错误率的说明（设**C=2**，**D=1**）

- 再假设t为唯一分界点（**C=2**，**D=1**）

- 做 ω_1 判别时的错误概率

$$P(e_{12}) = \int_{-\infty}^t P(\omega_2) p(x | \omega_2) dx$$

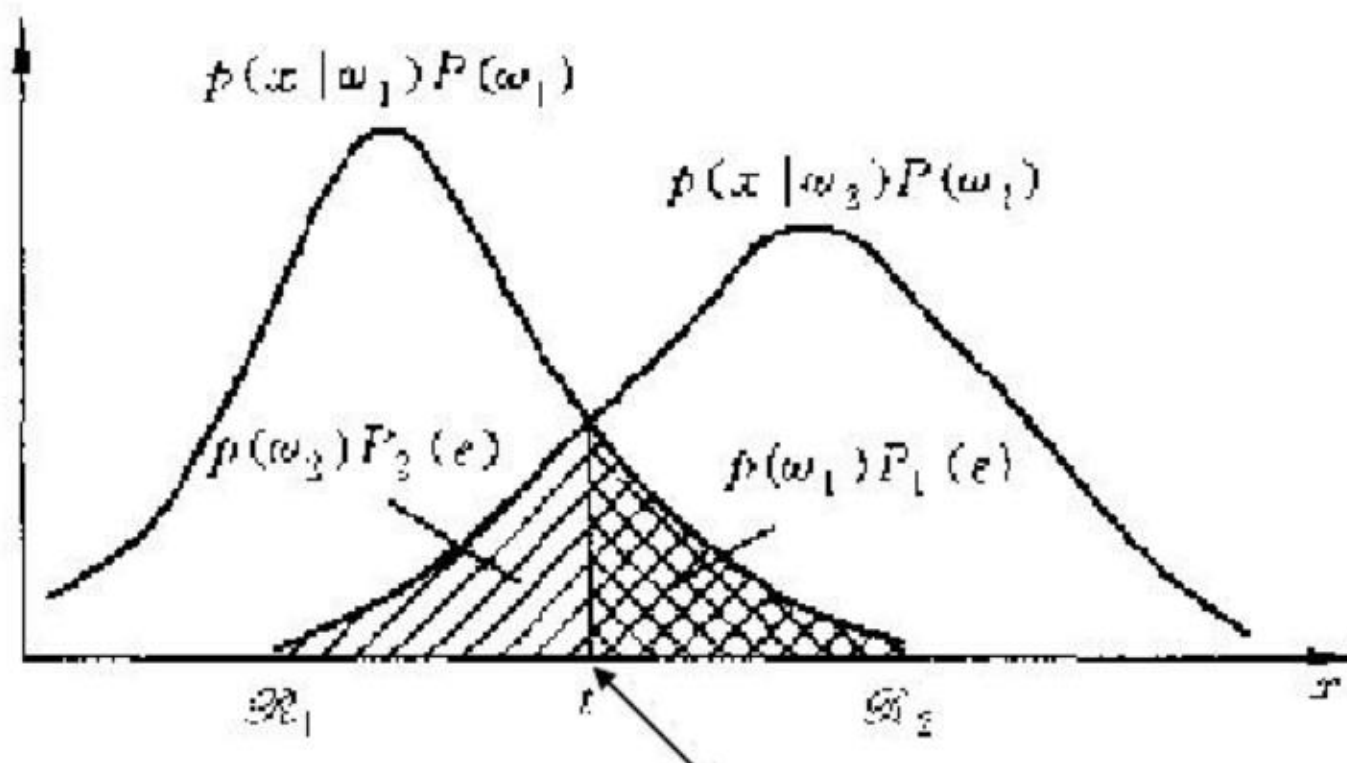
- 做 ω_2 判别时的错误概率

$$P(e_{21}) = \int_t^{\infty} P(\omega_1) p(x | \omega_1) dx$$

- 总错误概率

$$P(e) = \int_{-\infty}^t P(\omega_2) p(x | \omega_2) dx + \int_t^{\infty} P(\omega_1) p(x | \omega_1) dx$$

4.2.5 最小错误率的说明



$$P_1(e) = \int_{R_2} p(\mathbf{x} | \omega_1) d\mathbf{x}$$

$$P_2(e) = \int_{R_1} p(\mathbf{x} | \omega_2) d\mathbf{x}$$

4.2.5 最小错误率的说明

- 推广到多类（任意**C**类）
 - 已知条件和问题
 - **C**类**D**维问题： ω_1 , ω_2 , \dots , ω_C
 - 先验概率： $P(\omega_1)$, $P(\omega_2)$, \dots , $P(\omega_C)$
 - 类条件概率密度函数： $p(x | \omega_1)$, $p(x | \omega_2)$, \dots , $p(x | \omega_C)$
 - 发生了一个随机事件，其观察值为：特征向量 x
 - 求最小错误率分类器

4.2.5 最小错误率的说明

- 推广到多类（任意**C**类）
 - 判别函数
 - $P(\omega_i|x) = p(x|\omega_i) P(\omega_i) / p(x) \quad i = 1, \dots, C$
 - 决策规则
 - 对于未知样本 x ，若 $P(\omega_j|x) = \max P(\omega_i|x)$ ，则 $x \in \omega_j$

4.2.6 特点

- 最小错误率**Bayes**决策的特点
 - 已知条件多——各类概率分布
 - 最小错误率——概率意义上最优
 - 非线性分类器
 - 设计过程复杂

4.3 最小风险Bayes决策

- 4.3.1 问题的提出
- 4.3.2 决策规则
- 4.3.3 其它说明
- 4.3.4 特点

4.3.1 问题的提出

- 最小错误率**Bayes**决策
 - 最小错误率——概率意义上最优
 - 工程上是否最优？
- 错误分类的结果、代价或风险会是怎样的？
 - 考虑癌细胞图像识别的例子
- 出错的可能情况
 - 正常细胞 ω_1 错分为异常 ω_2
 - 异常细胞 ω_2 错分为正常 ω_1

4.3.1 问题的提出

- 区别状态和决策概念

- 状态：识别的目的是分类，把样本归类于其可能的自然状态（即类别）之一，将这种自然状态简称为状态，记为 ω
- 状态空间：所有可能的状态的集合构成状态空间，记为 Ω

4.3.1 问题的提出

- 区别状态和决策概念

- 决策：把样本归类于某个状态，或不能进行归类，都是决策，记为 α
- 决策空间：所有可能的决策（包括拒绝决策）的集合构成决策空间，记为 A

4.3.1 问题的提出

- 已知条件（类别**C = 2**，决策**A = 2**）
 - ω_1 、 $P(\omega_1)$ 、 $p(x | \omega_1)$ 和 ω_2 、 $P(\omega_2)$ 、 $p(x | \omega_2)$
 - $\Omega = \{\omega_1, \omega_2\}$
 - $A = \{\alpha_1, \alpha_2\}$
 - 定义损失函数 $\lambda(\alpha_i, \omega_j)$ ，简记 λ_{ij}
 - 发生了一个随机事件，其观察值为特征向量 x
- 求最小风险分类器

4.3.2 决策规则

- 判别函数

- $P(\omega_j | x) = p(x | \omega_j) P(\omega_j) / p(x) \quad j = 1, 2$

- $R(\alpha_i | x) = E[\lambda(\alpha_i, \omega_j) | x] = \lambda_{i1} P(\omega_1 | x) + \lambda_{i2} P(\omega_2 | x) \quad i = 1, 2$

- 决策规则

- 对于未知样本 x ，若 $R(\alpha_k | x) = \min R(\alpha_i | x)$ ，则 $x \in \omega_k$ ，即决策 α_k

4.3.2 决策规则

- 推广到任意情况（**C**个类别，**A**个决策）
 - ω_j 、 $P(\omega_j)$ 、 $p(x|\omega_j)$ $j = 1, \dots, C$
 - $\Omega = \{\omega_1, \omega_2, \dots, \omega_C\}$
 - $A = \{\alpha_1, \alpha_2, \dots, \alpha_A\}$
 - 定义损失函数 $\lambda(\alpha_i, \omega_j)$ ，简记 λ_{ij}
 - 发生了一个随机事件，其观察值为特征向量 x
 - 求最小风险分类器

4.3.2 决策规则

- 判别函数

- $P(\omega_j | x) = p(x | \omega_j) P(\omega_j) / p(x) \quad j = 1, 2, \dots, C$

- $R(\alpha_i | x) = E[\lambda(\alpha_i, \omega_j) | x] = \sum_{j=1} \lambda(\alpha_i, \omega_j) P(\omega_j | x) \quad i = 1, 2, \dots, A$

- 决策规则

- 对于未知样本 x ，若 $R(\alpha_k | x) = \min R(\alpha_i | x)$ ，则 $x \in \omega_k$ ，即决策 α_k

4.3.3 其它说明

- 最小风险与最小错误率的关系

- 0-1损失函数

$$\begin{aligned}\lambda(\alpha_i, \omega_j) &= 0 & i=j \\ \lambda(\alpha_i, \omega_j) &= 1 & i \neq j\end{aligned}$$

$$R(\alpha_i, x) = E[\lambda(\alpha_i, \omega_j)] = \sum_{j=1} \lambda(\alpha_i, \omega_j) P(\omega_j | x)$$

$$R(\alpha_i, x) = \sum_{\substack{j=1 \\ j \neq i}} P(\omega_j | x) = 1 - P(\omega_i | x)$$

4.3.3 其它说明

- 最小风险与最小错误率的关系

$$R(\alpha_i, x) = \sum_{\substack{j=1 \\ j \neq i}}^L P(\omega_j | x) = 1 - P(\omega_i | x)$$

- 对于未知样本 \mathbf{x} ，若 $R(\alpha_k | \mathbf{x}) = \min R(\alpha_i | \mathbf{x})$ ，则 $\mathbf{x} \in \omega_k$
- 对于未知样本 \mathbf{x} ，若 $P(\omega_k | \mathbf{x}) = \max P(\omega_i | \mathbf{x})$ ，则 $\mathbf{x} \in \omega_k$
- 结论：最小错误率**Bayes**决策，等价于**0-1**损失函数的最小风险**Bayes**决策

4.3.4 特点

- 最小风险**Bayes**决策的特点
 - 已知条件多——各类概率分布及风险系数
 - 最小错误风险——概率意义上最优
 - 非线性分类器
 - 设计过程复杂

例题

- 1、已知
 - 甲类: $P(\omega_1) = 0.7$ 和类条件概率密度函数 $p(x|\omega_1)$
 - 乙类: $P(\omega_2) = 0.3$ 和类条件概率密度函数 $p(x|\omega_2)$
 - 今有待分类样本特征观察值 $x = 10$, 且由函数曲线查得 $p(10|\omega_1) = 0.2$, $p(10|\omega_2) = 0.5$
 - (1)试用最小错误率Bayes决策对样本 $x = 10$ 进行分类
 - (2)试用最小风险Bayes决策对该样本进行分类, 设 $\lambda_{11}=\lambda_{22}=0$, $\lambda_{12}=2$, $\lambda_{21}=1$

4.4 最小最大决策

- 4.4.1 问题的提出
- 4.4.2 期望损失
- 4.4.3 最小最大风险
- 4.4.4 特点

4.4.1 问题的提出

- 已知条件和问题（**C = 2**情况）
 - 先验概率：考虑 $P(\omega_1)$ 和 $P(\omega_2)$ 未知或不确定的情况
 - 此时绝对意义的最小风险不存在
 - 如何求Bayes分类器
- 思路
 - 假设 $P(\omega_1)$ 和 $P(\omega_2)$ 确定
 - 设计一系列最小风险Bayes分类器
 - 取其中最大风险为最小的一个来用
 - 目的是控制最大风险

4.4.1 问题的提出

- 问题（类别**C = 2**，决策**A = 2**）
 - ω_1 、 $p(x | \omega_1)$ 和 ω_2 、 $p(x | \omega_2)$
 - $\Omega = \{\omega_1, \omega_2\}$
 - $A = \{\alpha_1, \alpha_2\}$
 - 损失函数 $\lambda(\alpha_i, \omega_j)$ ，简记 λ_{ij}
 - 发生了一个随机事件，其观察值为特征向量 x
- 求最小最大风险分类器

4.4.2 期望损失

- 期望损失的推导（回顾最小错误率证明）

$$P(e) = \int_{\mathfrak{R}_1} P(\omega_2) p(x | \omega_2) dx + \int_{\mathfrak{R}_2} P(\omega_1) p(x | \omega_1) dx$$

$$P(e_{12}) = \int_{\mathfrak{R}_1} P(\omega_2) p(x | \omega_2) dx$$

$$P(e_{21}) = \int_{\mathfrak{R}_2} P(\omega_1) p(x | \omega_1) dx$$

4.4.2 期望损失

- 期望损失的推导

联合或条件的
形式均可

$$\begin{aligned} R &= \int R(a(x)|x)p(x)dx = \int_{\mathcal{X}_1} R(a_1|x)p(x)dx + \int_{\mathcal{X}_2} R(a_2|x)p(x)dx \\ &= \int_{\mathcal{X}_1} [\lambda_{11}P(\omega_1)p(x|\omega_1) + \lambda_{12}P(\omega_2)p(x|\omega_2)]dx + \int_{\mathcal{X}_2} [\lambda_{21}P(\omega_1)p(x|\omega_1) \\ &\quad + \lambda_{22}P(\omega_2)p(x|\omega_2)]dx \end{aligned}$$

$P(\omega_2) = 1 - P(\omega_1)$, 代入

4.4.2 期望损失

- 期望损失的推导

$$R = \lambda_{22} + (\lambda_{12} - \lambda_{21}) \int_{\mathcal{R}_1} p(\mathbf{x} | \omega_2) d\mathbf{x} + P(\omega_1) [(\lambda_{11} - \lambda_{22}) \\ + (\lambda_{21} - \lambda_{11}) \int_{\mathcal{R}_2} p(\mathbf{x} | \omega_1) d\mathbf{x} - (\lambda_{12} - \lambda_{22}) \int_{\mathcal{R}_1} p(\mathbf{x} | \omega_2) d\mathbf{x}]$$

$$a = \lambda_{22} + (\lambda_{12} - \lambda_{21}) \int_{\mathcal{R}_1} p(\mathbf{x} | \omega_2) d\mathbf{x}$$

$$b = (\lambda_{11} - \lambda_{22}) + (\lambda_{21} - \lambda_{11}) \int_{\mathcal{R}_2} p(\mathbf{x} | \omega_1) d\mathbf{x} - (\lambda_{12} - \lambda_{22}) \int_{\mathcal{R}_1} p(\mathbf{x} | \omega_2) d\mathbf{x}$$

推导思路：替换后改成只有 $P(\omega_2)$ 有关的形式

4.4.3 最小最大风险

- 最小最大风险的解

- 任给 $P(\omega_1)^*$ 的值, $P(\omega_2)^* = 1 - P(\omega_1)^*$
 - ↓
- 求对应的最小风险Bayes分类器, 设决策域为 R_1 和 R_2
 - ↓
- 计算对应的 a^* 和 b^* , 得到 $R^* = a^* - b^* P(\omega_1)^*$ (注意! $R_1 R_2$ 随 $P(w)$ 变, 即决策会根据先验改变。故 R 和 $P(w)$ 不是线性)
 - ↓
- 重复上述计算步骤, 得到一系列最小风险Bayes分类器, 并 可得 $P(\omega_1)^* \text{——} R^*$ 关系曲线
 - ↓
- 比较所有最大风险, 取其中最小的一个, 作为最终的分类器

4.4.4 最小最大决策

- 最小最大决策的特点
 - 已知条件多——各类概率分布及风险系数
 - 最小最大风险——概率意义上最优
 - 非线性分类器
 - 设计过程很复杂

4.5 Bayes分类器设计

- 4.5.1 原理设计
- 4.5.2 Bayes决策面
- 4.5.3 错误率估计
- 4.5.4 其它

4.5.1 原理设计

- 两类情况（**C=2**）设计方法一

- 定义判别函数 $g(x)$

- ① $g(x) = P(\omega_1 | x) - P(\omega_2 | x)$

- ② $g(x) = p(x | \omega_1)P(\omega_1) - p(x | \omega_2)P(\omega_2)$

- ③ $g(x) = \ln \frac{p(x | \omega_1)}{p(x | \omega_2)} + \ln \frac{P(\omega_1)}{P(\omega_2)}$

- 决策面H方程

- $g(x) = 0$

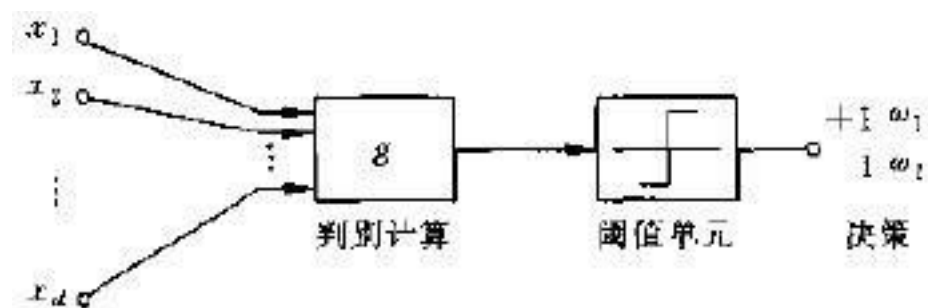
4.5.1 原理设计

- 两类情况 (**C=2**) 设计方法一

- 决策规则

- 对于未知样本 x , 若 $g(x) > 0$, 则 x 决策为 ω_1 类
- 若 $g(x) < 0$, 则 x 决策为 ω_2 类

- 原理图



4.5.1 原理设计

- 两类情况（**C=2**）设计方法二

- 定义判别函数 $g_1(x)$ 和 $g_2(x)$

- ① $g_i(x) = P(\omega_i | x)$

- ② $g_i(x) = p(x | \omega_i) P(\omega_i)$

- ③ $g_i(x) = \ln p(x | \omega_i) + \ln P(\omega_i)$

- 决策面H方程

- $g_1(x) = g_2(x)$

4.5.1 原理设计

- 两类情况（**C=2**）设计方法二
 - 决策规则
 - 对于未知样本 x ，若 $g_1(x) > g_2(x)$ ，则 x 决策为 ω_1 类
 - 若 $g_1(x) < g_2(x)$ ，则 x 决策为 ω_2 类
 - 原理图

4.5.1 原理设计

- 多类情况（**C**任意）设计方法
 - 定义判别函数 $g_i(\mathbf{x})$ $i = 1, 2, \dots, C$
 - ① $g_i(\mathbf{x}) = P(\omega_i | \mathbf{x})$
 - ② $g_i(\mathbf{x}) = p(\mathbf{x} | \omega_i) P(\omega_i)$
 - ③ $g_i(\mathbf{x}) = \ln p(\mathbf{x} | \omega_i) + \ln P(\omega_i)$
 - 决策面H方程
 - $g_i(\mathbf{x}) = g_j(\mathbf{x})$ $i, j = 1, 2, \dots, C$

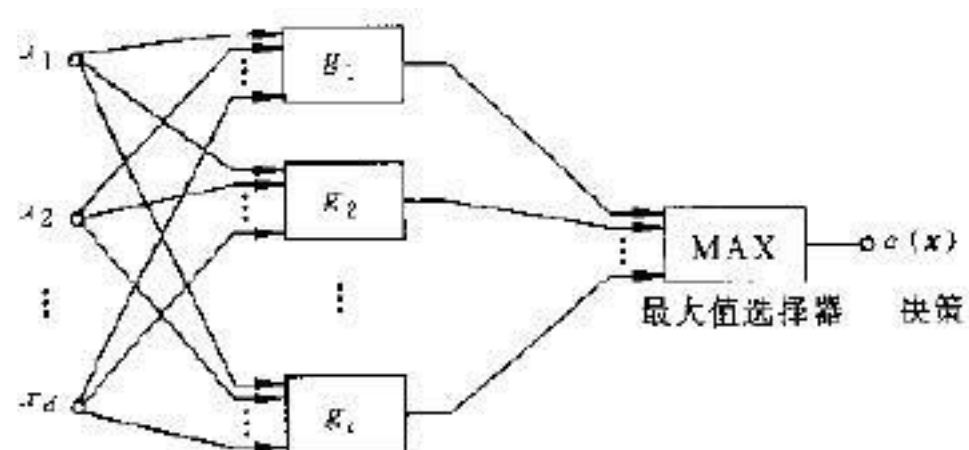
4.5.1 原理设计

- 多类情况（**C**任意）设计方法

- 决策规则

- 对于未知样本 x ，若 $g_j(x) = \text{MAX } g_i(x)$ ，则 x 决策为 ω_j 类

- 原理图



4.5.2 正态分布决策面

- 单变量正态分布 / 一元正态分布

$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right\}$$

4.5.2 正态分布决策面

- 多变量正态分布 / 多元正态分布

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \triangleq \frac{1}{(2\pi)^{D/2}|\boldsymbol{\Sigma}|^{1/2}} \exp \left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \right]$$

$$\boldsymbol{\mu} = E\{\mathbf{x}\}$$

$$\boldsymbol{\Sigma} = E\{(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T\}$$

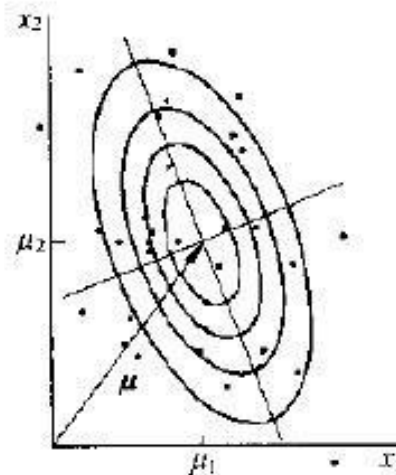
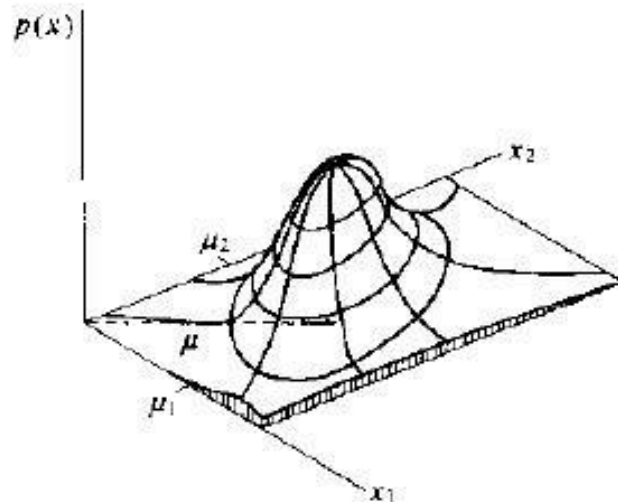
4.5.2 正态分布决策面

- 多元正态分布的性质

- 均值与方差
- 等概率密度点的轨迹

正态分布的等密度点的轨迹为超椭球面

$$(x - \mu)^T \Sigma^{-1} (x - \mu) = \text{常数}$$



$$\gamma^2 = (x - \mu)^T \Sigma^{-1} (x - \mu)$$

4.5.2 正态分布决策面

- 多元正态分布的性质

- 不相关性 & 独立性: 不相关性 = 独立性

$$\begin{aligned}\sigma_{ij}^2 &= E[(x_i - \mu_i)(x_j - \mu_j)], \quad i, j = 1, 2, \dots, d, i \neq j \\ &= E(x_i - \mu_i) \cdot E(x_j - \mu_j) \\ &= 0\end{aligned}$$

$$\Sigma = \begin{bmatrix} \sigma_{11}^2 & \cdots & 0 \\ \vdots & \ddots & \\ 0 & \cdots & \sigma_{dd}^2 \end{bmatrix}$$

$$\Sigma^{-1} = \begin{bmatrix} 1/\sigma_{11}^2 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & 1/\sigma_{dd}^2 \end{bmatrix}$$

4.5.2 正态分布决策面

- 多元正态分布的性质

- 边缘分布: 正态分布
- 条件分布: 正态分布
- 线性组合: 正态分布

4.5.2 正态分布决策面

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu) \right\}$$

- 正态分布时分类器的决策面方程

- 判别函数

③ $g_i(\mathbf{x}) = \ln p(\mathbf{x} | \omega_i) + \ln P(\omega_i)$

$$g_i(\mathbf{x}) = -\frac{1}{2} (\mathbf{x} - \mu_i)^T \Sigma_i^{-1} (\mathbf{x} - \mu_i) - \frac{d}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma_i| + \ln P(\omega_i)$$

- 决策面方程

$$g_i(\mathbf{x}) = g_j(\mathbf{x})$$

$$-\frac{1}{2} [(\mathbf{x} - \mu_i)^T \Sigma_i^{-1} (\mathbf{x} - \mu_i) - (\mathbf{x} - \mu_j)^T \Sigma_j^{-1} (\mathbf{x} - \mu_j)] - \frac{1}{2} \ln \frac{|\Sigma_i|}{|\Sigma_j|} + \ln \frac{P(\omega_i)}{P(\omega_j)} = 0$$

4.5.2 正态分布决策面

- 正态分布时分类器的决策面方程

- 二次型判别函数

展开

$$\begin{aligned}g_i(\mathbf{x}) &= -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}_i^{-1}(\mathbf{x} - \boldsymbol{\mu}_i) - \frac{1}{2} \ln |\boldsymbol{\Sigma}_i| + \ln P(\omega_i) \\&= \mathbf{x}^T \mathbf{W}_i \mathbf{x} + \mathbf{w}_i^T \mathbf{x} + w_{i0}\end{aligned}$$

$$\mathbf{W}_i = -\frac{1}{2} \boldsymbol{\Sigma}_i^{-1} \quad (d \times d \text{ 矩阵})$$

$$\mathbf{w}_i = \boldsymbol{\Sigma}_i^{-1} \boldsymbol{\mu}_i \quad (d \text{ 维列向量})$$

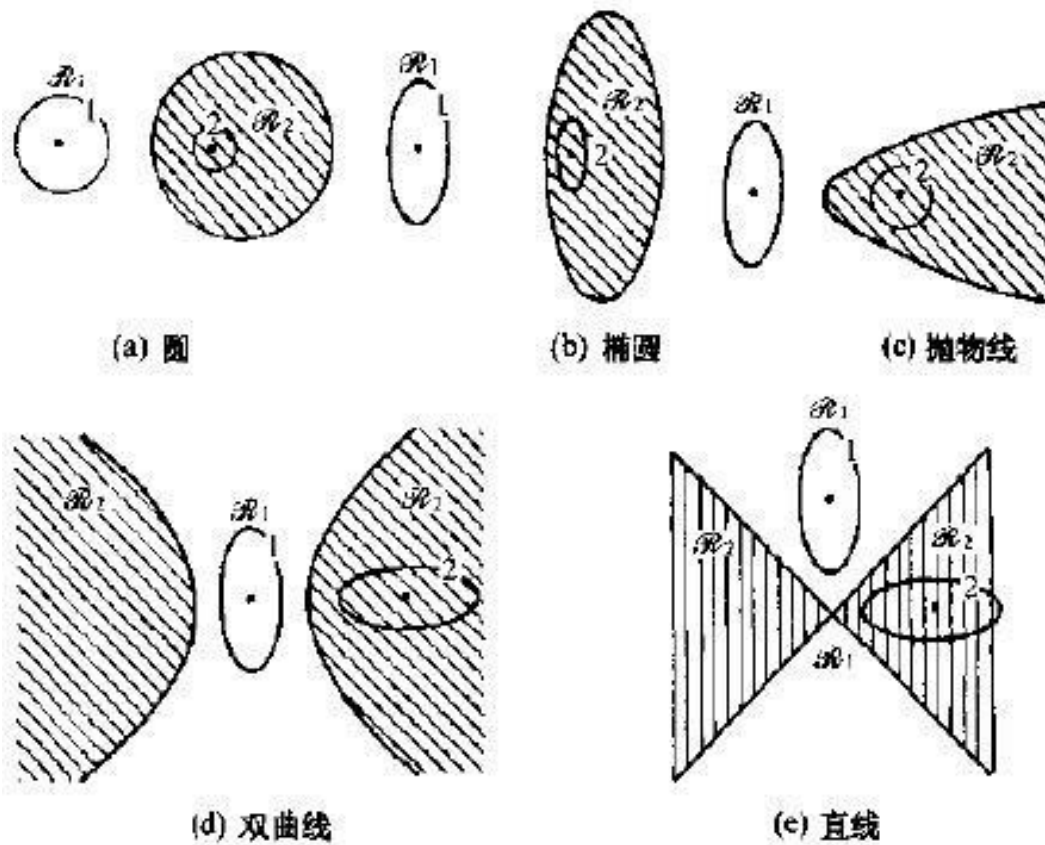
$$w_{i0} = -\frac{1}{2} \boldsymbol{\mu}_i^T \boldsymbol{\Sigma}_i^{-1} \boldsymbol{\mu}_i - \frac{1}{2} \ln |\boldsymbol{\Sigma}_i| + \ln P(\omega_i)$$

- 决策面方程

相减

4.5.2 正态分布决策面

- 决策面示例





THE END !