

模式识别

干晓蓉 编著

云南出版集团公司
云南人民出版社

模式识别

于晓蓉 编著

云南出版集团公司
云南人民出版社

图书在版编目 (CIP) 数据

模式识别 / 干晓蓉编著. —昆明: 云南人民出版社,
2006.9

ISBN 7-222-04857-X

I. 模... II. 干... III. 模式识别 IV. 0235

中国版本图书馆 CIP 数据核字 (2006) 第 118299 号

责任编辑: 朱海涛 王绍来
装帧设计: 黄麟 孙勇
责任印制: 马跃武

书名	模式识别
作者	干晓蓉 编著
出版	云南出版集团公司 云南人民出版社
发行	云南人民出版社
社址	昆明市环城西路 609 号
邮编	650034
网址	www.ynpph.com.cn
E-mail	rmszbs@public.km.yn.cn
开本	787×1092 1/16
印张	11.75
字数	250 千字
版次	2006 年 9 月第 1 版第 1 次印刷
印刷	昆明宏宝兴印务有限公司
书号	ISBN7-222-04857-X
定价	26.00 元

尊敬的读者: 若你购买的我社图书存在印装质量问题, 请与我社发行部联系调换。
电话: (0871) 4194864 4191604 4107628 (邮购)

序 言

模式识别 (Pattern Recognition) 是指对表征事物或现象的各种形式的 (数值的、文字和逻辑关系的) 信息进行处理和分析, 以对事物或现象进行描述、辨认、分类和解释的过程。作为一门学科, 它研究用机器完成自动识别事物的工作, 是信息科学和人工智能的重要组成部分。

模式识别的学习和研究主要集中于两个方面: 一是研究生物体 (包括人类) 是如何感知对象的, 属于认识科学范畴, 这是生物学家、心理学家、生物学家和神经生理学家们的研究对象。二是在给定任务下, 如何通过建立数学模型, 用高速计算机实现模式识别的理论和方法, 这方面是数学家、信息理论专家和计算机工作者的工作和研究方向。

本书是编著者为昆明理工大学理工科相关专业的研究生和高年级本科生讲授模式识别的理论和方法课程的总结。主要讨论统计识别理论和方法, 包括 Bayes 决策理论、概率密度估计、线性判别函数、无监督学习和聚类、特征选择与提取等内容。并讨论了模糊模式识别、人工神经网络等课题。每章末都附有习题, 并在篇末增加了线性代数和概率统计的基本知识作为附录。

著名数学家和社会活动家、北京大学教授丁石孙说过这样一段话: “学习数学当然要学习一些理论, 学习一些定理与概念, 也要学习一些解题的技巧。但更重要的是学到数学的思想方法, 用以解决数学和数学以外的问题, 特别是要学会用数学来解决许多非数学问题。” 模式识别其实就是用数学方法解决实际问题的理论与应用相结合的新学科, 学习这门课的时候, 如果不通过实际例子来理解和认识各种方法, 其学习效果是不好的。本书编著者充分理解上述丁教授的一段话, 在讲述每种识别方法时, 都首先理解所涉及的概念, 并配有生动直观的例子, 以及充分应用图解来阐明理论和方法, 取材是独具匠心的。

笔者热忱推荐本书的正式出版发行。它将作为本课程教材建设的一项新成果而使广大读者受益, 谨此为序。

李继彬

2006 年 8 月

目 录

第一章 绪 论	(1)
1.1 模式识别的概念	(1)
1.2 模式识别的方法	(2)
第二章 贝叶斯决策理论	(4)
2.1 最小错误贝叶斯决策规则	(4)
2.2 最小风险贝叶斯决策规则	(9)
2.3 聂曼—皮尔逊决策规则	(12)
2.4 最小最大决策	(15)
2.5 正态分布的统计决策	(16)
2.5.1 多元正态分布	(17)
2.5.2 正态分布模式的贝叶斯分类器	(18)
2.5.3 正态分布且等协方差时的错误率计算	(20)
习题	(22)
第三章 概率密度估计	(25)
(一) 参数估计	(25)
3.1 均值向量和协方差矩阵的矩法估计	(26)
3.2 最大似然估计	(27)
3.3 贝叶斯估计	(30)
3.3.1 参数的贝叶斯估计	(30)
3.3.2 概率密度的贝叶斯估计	(31)
3.3.3 递推的贝叶斯估计	(32)
(二) 非参数估计	(37)
3.4 直方图法	(37)
3.5 k 近邻法	(38)
3.5.1 k 近邻法估计密度	(38)
3.5.2 k 近邻决策规则	(39)
3.5.3 最近邻决策的错误率	(40)
3.6 用基函数展开法	(40)

3.7 Parzen 窗法	(43)
3.7.1 Parzen 窗估计法	(43)
3.7.2 窗函数的选择	(44)
3.7.3 估计量 $\hat{p}_n(\boldsymbol{x})$ 为密度函数的条件	(44)
3.7.4 窗宽 h_n 对估计量 $\hat{p}_n(\boldsymbol{x})$ 的影响	(45)
3.7.5 估计量 $\hat{p}(\boldsymbol{x})$ 的统计性质	(45)
习题	(47)
第四章 线性判别函数	(50)
4.1 线性判别函数的基本概念	(50)
4.1.1 线性判别函数	(50)
4.1.2 两类问题的决策规则和决策面	(50)
4.1.3 广义线性判别函数	(52)
4.1.4 设计线性分类器的主要步骤	(53)
4.2 Fisher 线性判别	(54)
4.3 感知器准则函数	(57)
4.3.1 几个基本概念	(57)
4.3.2 感知器准则函数及其梯度下降算法	(59)
4.4 最小平方误差准则函数	(63)
4.4.1 平方误差准则函数及其伪逆解	(63)
4.4.2 与 Fisher 线性判别的关系	(65)
4.4.3 对贝叶斯判别函数的渐近逼近	(67)
4.4.4 MSE 准则函数的梯度下降算法	(68)
4.5 支持向量机	(69)
4.5.1 样本线性可分	(69)
4.5.2 样本线性不可分	(74)
4.6 多类问题	(75)
习题	(77)
第五章 无监督学习和聚类	(79)
5.1 无监督的参数估计	(79)
5.1.1 混合密度及其可辨识性	(79)
5.1.2 最大似然估计	(80)
5.1.3 对混合正态密度的应用	(83)

5.2	聚类分析的基本概念	(85)
5.2.1	模式间的相似测度	(85)
5.2.2	聚类准则函数	(88)
5.3	聚类算法	(91)
5.3.1	层次聚类法	(92)
5.3.2	c -均值算法	(93)
	习题	(97)
第六章	特征选择与提取	(100)
6.1	引言	(100)
6.2	基于几何距离的特征提取	(100)
6.2.1	基于距离的类别可分离性判据	(100)
6.2.2	按基于几何距离判据的特征提取方法	(102)
6.3	基于概率分布的特征提取	(105)
6.3.1	基于概率分布的可分性判据	(105)
6.3.2	用散度准则函数 J_D 的特征提取方法	(108)
6.4	基于 $K-L$ 变换的特征提取	(111)
6.4.1	$K-L$ 变换	(111)
6.4.2	基于 $K-L$ 变换的特征提取	(113)
6.5	特征选择	(118)
6.5.1	最优搜索算法	(119)
6.5.2	次优搜索算法	(120)
	习题	(122)
第七章	模糊模式识别	(124)
7.1	引言	(124)
7.2	模糊集合及其运算	(124)
7.2.1	普通集合及其特征函数	(124)
7.2.2	模糊集合的定义及其表示	(125)
7.2.3	模糊集合的基本运算	(127)
7.2.4	模糊集合的 λ 截集	(129)
7.3	模糊模式识别的一些方法	(130)
7.3.1	模糊模式识别的直接方法	(130)
7.3.2	基于择近原则的模糊模式识别	(132)

7.4	模糊关系	(135)
7.4.1	普通集合间的关系	(135)
7.4.2	模糊关系及模糊矩阵	(137)
7.5	模糊聚类分析	(141)
7.5.1	基于模糊等价关系的聚类方法	(141)
7.5.2	模糊 c -均值算法	(145)
	习题	(148)
第八章	人工神经网络	(150)
8.1	单层感知器网络	(150)
8.1.1	两类问题	(150)
8.1.2	多类情况	(152)
8.1.3	激励函数	(152)
8.2	前向多层神经网络	(153)
8.2.1	异或问题	(153)
8.2.2	前向多层神经网络的 BP 算法	(155)
8.2.3	关于 BP 算法的几点注解	(159)
8.3	广义线性分类器	(160)
8.3.1	广义线性分类器的结构	(160)
8.3.2	高阶神经网络	(161)
8.3.3	径向基函数网络	(163)
附录 A	线性代数	(166)
A.1	矩阵运算	(166)
A.2	向量的内积与外积	(168)
A.3	矩阵的本征值与本征向量	(168)
A.4	矩阵的导数	(170)
附录 B	多维随机变量	(173)
B.1	二维情形 (X, Y)	(173)
B.2	多维情形 $X = (X_1, \dots, X_d)$	(173)
B.3	数学期望、均值向量和协方差矩阵	(175)
	参考书目	(177)

第一章 绪 论

模式识别诞生于20世纪20年代,随着40年代计算机的出现,50年代人工智能的兴起,模式识别作为一个研究领域,迅速发展于20世纪60年代。它是一个多领域的交叉学科,涉及到统计学、工程学、人工智能、计算机科学、心理学和生理学等学科。它所研究的理论和方法在很多科学技术领域中得到了广泛的应用,推动了人工智能系统的发展,扩大了计算机应用的可能性。许多人为了解决实际问题进入该领域,其中包括字符自动识别、医疗诊断等经典问题和个人信用评分、商品销售分析、信用卡交易分析等关于数据挖掘的新问题。如此广泛的模式识别应用,吸引了众多的研究力量,产生出许多新方法,推动着该学科的进一步发展。

几十年来,模式识别研究取得了大量的成果,在很多地方得到了成功的应用。但是,由于模式识别涉及到很多复杂的问题,现有的理论和方法对于解决这些问题还有很多不足之处。因此模式识别仍然是一门发展中的新兴学科,新的理论和方法不断出现。同时,与其他学科相互结合相互渗透,不断推动模式识别向前发展。

什么是模式识别?模式识别是所有生物所具有的特性,然而不同的生物,其识别的方法不同。人能够通过视觉、声音或手迹轻而易举地辨别人脸、识别语音、阅读手写文字、从口袋里摸出钥匙,或者根据气味判断苹果是否成熟。狗可以闻到30码以外的生物的气味,这一点人做不到。然而多数的狗对照镜子却没什么深刻印象,因为它实际上并没有真正意识到镜子里的另一只狗。再如当微生物来到pH值不适合的环境中时就会逃走。然而物体的识别并不局限在生物的感官中,在一个谈话里,我们能突然分辨出多年前曾听到过的谈话。所有这些例子,就是模式识别(模式分类)。这种输入原始数据并根据其类别采取相应行为的能力,对于我们的生存至关重要。为了具有这种能力,在过去的几千万年里,我们进化出高度复杂的神经和认知系统。随着计算机的广泛应用,人们希望计算机也能有这种识别能力,试图设计和建造一台能够识别不同模式的机器的想法是很自然的。从自动语音识别到指纹识别、光学字符识别、DNA序列分析等等很多的应用,都清楚地表明一个可靠和准确的模式识别机器的巨大作用。

1.1 模式识别的概念

模式——对客体的定量的或结构的描述。

定量描述:例如为了识别细胞是正常细胞还是癌细胞,抓住细胞的两个特征,即圆形度 x_1 ,形心偏差度 x_2 ,组成描述细胞的特征向量 $x = (x_1, x_2)^T = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$,这就是细胞定量描述的模式。

结构描述:例如图 1.1 的图形,由两个圆弧段 a, c 及两直线段 b 组成,这四个基元组成符号串 $x = abcb$,这就是这图形结构描述的模式。

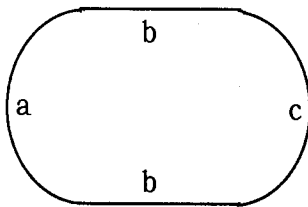


图 1.1

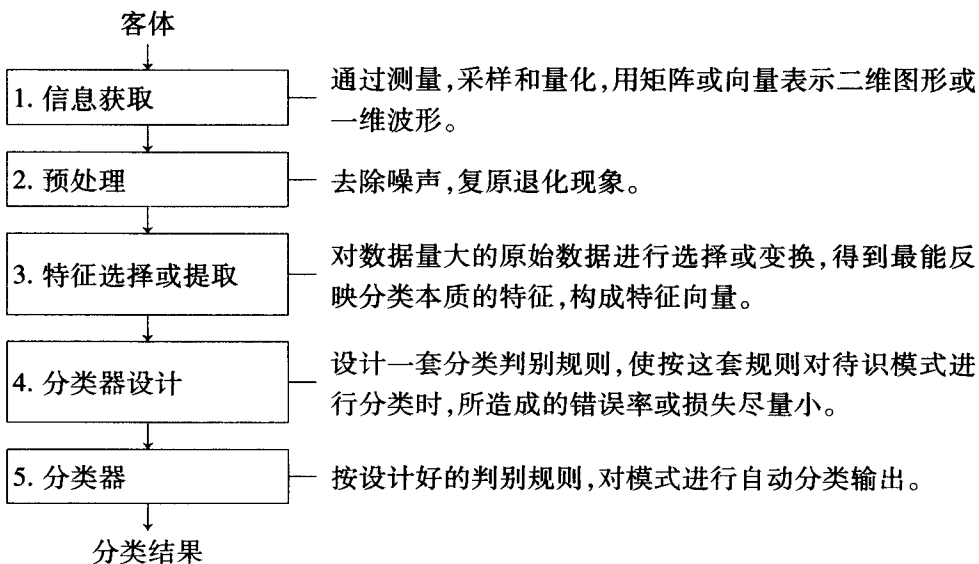
模式类——具有某些共同特性的模式的集合。例如所有细胞分为正常细胞与癌细胞两类,正常细胞的共同特性是圆形度大,细胞核偏离细胞中心小,癌细胞共同特性则与此相反。

模式识别——研究一些自动技术,使计算机能自动地把待识别的模式分到各自的模式类中去。

1.2 模式识别的方法

有两种基本的模式识别方法:统计模式识别方法及结构模式识别方法。统计模式识别方法用特征向量描述模式,结构模式识别方法用符号串描述模式。

模式识别系统由下面五部分组成:



本书只讨论后三部分的理论和方法。内容包括：贝叶斯决策理论、概率密度估计、线性判别函数、无监督学习和聚类、特征的选择与提取、模糊模式识别以及人工神经网络等内容。

第二章 贝叶斯决策理论

贝叶斯决策理论和方法是统计模式识别中的一个基本方法,采用这个方法进行模式分类时,要求满足以下两个条件:

(1) 决策分类的类别数是一定的。设有 c 类,记为 $\omega_1, \omega_2, \dots, \omega_c$;

(2) 各类别总体的概率分布是已知的。即各类别 ω_i 出现的先验概率 $p(\omega_i)$ 及类条件概率密度函数 $p(\mathbf{x} | \omega_i)$ 是已知的。

记号 $p(\omega_i | \mathbf{x})$ 表示模式 \mathbf{x} 出现的条件下, ω_i 类出现的概率,称为 ω_i 的后验概率,可理解为模式 \mathbf{x} 来自 ω_i 的概率。根据贝叶斯公式,有

$$p(\omega_i | \mathbf{x}) = \frac{p(\omega_i, \mathbf{x})}{p(\mathbf{x})} = \frac{p(\mathbf{x} | \omega_i)p(\omega_i)}{\sum_{j=1}^c p(\mathbf{x} | \omega_j)p(\omega_j)} \quad (2-1)$$

因此,后验概率 $p(\omega_i | \mathbf{x})$ 也是已知的。

2.1 最小错误贝叶斯决策规则

先讨论两类问题。例如,要进行细胞识别,设细胞模式为 $\mathbf{x} = (x_1, x_2)^T$,其中 x_1 为细胞的圆形度, x_2 为形心偏差度,要识别 \mathbf{x} 属正常类 ω_1 或属癌变类 ω_2 ? 因为我们假设模式 \mathbf{x} 来自两类的概率 $p(\omega_1 | \mathbf{x})$ 与 $p(\omega_2 | \mathbf{x})$ 是已知的,比较它们的大小,若 $p(\omega_1 | \mathbf{x})$ 大,则 \mathbf{x} 判属 ω_1 类,若 $p(\omega_2 | \mathbf{x})$ 大,则 \mathbf{x} 判属 ω_2 类,这就是最小错误的贝叶斯决策规则。

对一般两类问题,最小错误贝叶斯决策规则有以下几种等价形式:

(1) 后验概率形式

$$\text{若 } p(\omega_1 | \mathbf{x}) \geq p(\omega_2 | \mathbf{x}), \text{ 则 } \mathbf{x} \in \begin{cases} \omega_1 \\ \omega_2 \end{cases} \quad (2-2)$$

根据贝叶斯公式

$$p(\omega_i | \mathbf{x}) = \frac{p(\mathbf{x} | \omega_i)p(\omega_i)}{p(\mathbf{x})}, i = 1, 2$$

就有下面的等价形式

(2) 类条件密度形式

$$\text{若 } p(\mathbf{x} | \omega_1)p(\omega_1) \geq p(\mathbf{x} | \omega_2)p(\omega_2), \text{ 则 } \mathbf{x} \in \begin{cases} \omega_1 \\ \omega_2 \end{cases} \quad (2-3)$$

(3) 似然比形式

$$\text{若 } \frac{p(x | \omega_1)}{p(x | \omega_2)} \geq \frac{p(\omega_2)}{p(\omega_1)}, \text{ 则 } x \in \begin{cases} \omega_1 \\ \omega_2 \end{cases} \quad (2-4)$$

其中比值 $p(x | \omega_1)/p(x | \omega_2)$ 称为似然比, 而 $p(\omega_2)/p(\omega_1)$ 称为似然比的阈值。

有时取对数计算比较方便, 而 $\ln x$ 为单调增加函数, 对不等式取对数, 不改变左、右的大小关系, 因而又有

(4) 似然比对数形式

$$\text{若 } \ln \frac{p(x | \omega_1)}{p(x | \omega_2)} \geq \ln \frac{p(\omega_2)}{p(\omega_1)}, \text{ 则 } x \in \begin{cases} \omega_1 \\ \omega_2 \end{cases} \quad (2-5)$$

推广到 c 类情况, 贝叶斯决策规则为

(1) 后验概率形式

$$\text{若 } p(\omega_i | x) > p(\omega_j | x), j = 1, 2, \dots, c, j \neq i, \text{ 则 } x \in \omega_i \quad (2-6)$$

(2) 类条件概率密度形式

$$\text{若 } p(x | \omega_i)p(\omega_i) > p(x | \omega_j)p(\omega_j), j = 1, 2, \dots, c, j \neq i, \text{ 则 } x \in \omega_i \quad (2-7)$$

(3) 似然比形式

$$\text{若 } \frac{p(x | \omega_i)}{p(x | \omega_j)} > \frac{p(\omega_j)}{p(\omega_i)}, j = 1, 2, \dots, c, j \neq i, \text{ 则 } x \in \omega_i \quad (2-8)$$

(4) 似然比对数形式

$$\text{若 } \ln \frac{p(x | \omega_i)}{p(x | \omega_j)} > \ln \frac{p(\omega_j)}{p(\omega_i)}, j = 1, 2, \dots, c, j \neq i, \text{ 则 } x \in \omega_i \quad (2-9)$$

例 2.1 某医院研究癌症的诊断, 对一大批人打试验针作普查, 统计数字为:

- ① 每 1000 人有 5 个癌症病人;
- ② 每 100 个正常人有 1 人对试验反应阳性;
- ③ 每 100 个癌症病人有 95 人对试验反应阳性。

现某人试验结果为阳性, 诊断结论是什么?(正常人还是癌症病人?)

设 ω_1 类 —— 正常人, ω_2 类 —— 癌症病人

特征值 x 取为试验结果, 即 $x =$ 阳性或 $x =$ 阴性, 则

$$p(\omega_1) = 0.995, p(\omega_2) = 0.005, p(\text{阳} | \omega_1) = 0.01$$

$$p(\text{阴} | \omega_1) = 0.99, p(\text{阳} | \omega_2) = 0.95, p(\text{阴} | \omega_2) = 0.05$$

现设 $x =$ 阳性, 问 $x \in \omega_1$ 或 $x \in \omega_2$?

$$\text{因为 } p(\omega_1)p(x | \omega_1) = p(\omega_1)p(\text{阳} | \omega_1) = 0.995 \times 0.01 = 0.00995$$

$$p(\omega_2)p(x | \omega_2) = p(\omega_2)p(\text{阳} | \omega_2) = 0.005 \times 0.95 = 0.00475$$

$$p(\omega_1)p(x | \omega_1) > p(\omega_2)p(x | \omega_2)$$

所以, 根据最小错误率贝叶斯决策规则, $x \in \omega_1$ 类, 即某人属正常人。

判别函数是指模式 x 的某些函数,由它可以导出分类规则。例如,在两类问题中,可取两个判别函数

$$g_1(x) = p(x | \omega_1)p(\omega_1), \quad g_2(x) = p(x | \omega_2)p(\omega_2)$$

则贝叶斯分类规则为

$$\text{若 } g_1(x) \geq g_2(x), \text{ 则 } x \in \begin{cases} \omega_1 \\ \omega_2 \end{cases}$$

也可取判别函数为 $g(x) = g_1(x) - g_2(x)$, 则贝叶斯分类规则为

$$\text{若 } g(x) \geq 0, \text{ 则 } x \in \begin{cases} \omega_1 \\ \omega_2 \end{cases}$$

也可取似然比 $l(x) = \frac{p(x | \omega_1)}{p(x | \omega_2)}$ 作为判别函数,从而导出相应的判别规则。因此,判别函数不是唯一的。

在 c 类情况下,通常需要定义 c 个判别函数 $g_i(x)$, 而相应的分类规则为

若 $g_i(x) > g_j(x), j = 1, 2, \dots, c, j \neq i$, 则 $x \in \omega_i$

例如定义 $g_i(x) = p(x | \omega_i)p(\omega_i)$, 由此可导出贝叶斯决策规则。

模式向量 x 所有可能的取值形成模式空间 Ω , 决策规则将 Ω 分成 c 个决策域 $\Omega_1, \Omega_2, \dots, \Omega_c$, 在 Ω_i 中的模式, 决策规则都将其归于 ω_i 类。决策域的边界称为决策边界, 若 Ω_i 与 Ω_j 是相邻的决策域, 则决策边界的方程为 $g_i(x) = g_j(x)$ 。

模式 $x = (x_1)$ 为一维时, 决策边界为点;

模式 $x = (x_1, x_2)^T$ 为二维时, 决策边界为曲线;

模式 $x = (x_1, x_2, x_3)^T$ 为三维时, 决策边界为空间曲面;

模式 $x = (x_1, \dots, x_d)^T$ 为 d 维时, 决策边界为 d 维空间超曲面。

图 2.1 标明了模式为一维时, 两类问题的决策域 Ω_1, Ω_2 及决策边界 t 。

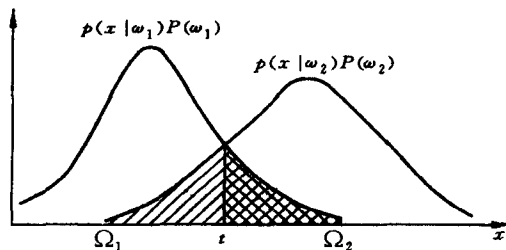


图 2.1

分类器可看作是由硬件和软件组成的一个“机器”, 它的功能是先计算出判别函数的值, 再根据决策规则从中选出一个类作为决策的结果。应用贝叶斯决策规则

对模式进行分类的分类器,称贝叶斯分类器。

从以下的讨论可以看出,前面所述的贝叶斯决策规则确实使分类错误最小。由全概率公式, c 类问题分类错误概率 $p(e)$ 可用下式表示:

$$p(e) = \sum_{i=1}^c p(e | \omega_i) p(\omega_i) \quad (2-10)$$

其中 $p(e | \omega_i)$ 是 ω_i 类的样本被错分的概率,它应该等于 ω_i 类的样本 \mathbf{x} 落入 Ω_i 的余集 $\Omega - \Omega_i$ 中的概率,即

$$p(e | \omega_i) = \int_{\Omega - \Omega_i} p(\mathbf{x} | \omega_i) d\mathbf{x}$$

于是有

$$\begin{aligned} p(e) &= \sum_{i=1}^c p(\omega_i) \int_{\Omega - \Omega_i} p(\mathbf{x} | \omega_i) d\mathbf{x} \quad (2-11) \\ &= \sum_{i=1}^c p(\omega_i) \left(1 - \int_{\Omega_i} p(\mathbf{x} | \omega_i) d\mathbf{x} \right) \\ &= 1 - \sum_{i=1}^c \int_{\Omega_i} p(\mathbf{x} | \omega_i) p(\omega_i) d\mathbf{x} \end{aligned}$$

要使错分概率 $p(e)$ 最小,等价于正分概率

$$\sum_{i=1}^c \int_{\Omega_i} p(\mathbf{x} | \omega_i) p(\omega_i) d\mathbf{x}$$

最大,也就是要选择积分区域 Ω_i ,使上式中的积分值最大,而贝叶斯决策规则正是满足

$$\mathbf{x} \in \Omega_i \text{ 时, } p(\mathbf{x} | \omega_i) p(\omega_i) = \max_{1 \leq j \leq c} p(\mathbf{x} | \omega_j) p(\omega_j)$$

因此,贝叶斯决策规则就是使分类错误最小的决策规则,此时正分概率可写成

$$a = \int_{\Omega} \max_j p(\mathbf{x} | \omega_j) p(\omega_j) d\mathbf{x} \quad (2-12)$$

错分概率为

$$p(e) = 1 - a = \int_{\Omega} [1 - \max_j p(\omega_j | \mathbf{x})] p(\mathbf{x}) d\mathbf{x} \quad (2-13)$$

对于两类问题,(2-11)式成为

$$p(e) = p(\omega_1) \int_{\Omega_2} p(\mathbf{x} | \omega_1) d\mathbf{x} + p(\omega_2) \int_{\Omega_1} p(\mathbf{x} | \omega_2) d\mathbf{x} \quad (2-14)$$

$$= p(\omega_1) p_1(e) + p(\omega_2) p_2(e) \quad (2-15)$$

(2-15)式右边第一项是 ω_1 类的样本错分的概率,等于图2.1中右边纹线所示的面积,第二项是 ω_2 类的样本错分的概率,等于图2.1中左边斜线所示的面积。

对于在决策边界面上及其附近样本的分类,容易导致分类错误,因此拒绝对某些样本作出决策可以使错误率降低,但这样做也可能使得本来正确分类的样本也

遭到拒绝。这里就要讨论错误率与拒绝率的权衡问题。

首先,将模式空间 Ω 分成两个互补的区域如下

$$\text{拒绝域 } R = \{x \mid \max_i p(\omega_i | x) < 1 - t\} \quad (2-16)$$

$$\text{接受域 } A = \{x \mid \max_i p(\omega_i | x) \geq 1 - t\} \quad (2-17)$$

其中 t 为阈值。图 2.2 是某个两类问题的后验概率曲线,图中标出了拒绝域,接受域和阈值。

因为 $1 = \sum_{i=1}^c p(\omega_i | x) \leq c \cdot \max_i p(\omega_i | x)$, $\max_i p(\omega_i | x) \geq \frac{1}{c}$,所以要使拒绝域非空,必须取阈值 t 满足

$$t < 1 - \max_i p(\omega_i | x) \leq 1 - \frac{1}{c} = \frac{c-1}{c} \quad (2-18)$$

例如图 2.2 中的两类问题,应使 $t < \frac{1}{2}$ 或 $1 - t > \frac{1}{2}$ 。

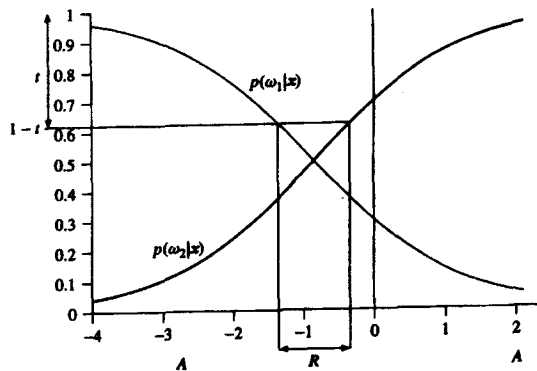


图 2.2

确定了接受域 A 及拒绝域 R 后,对位于接受域 A 内的模式 x ,可用基于最小错误的贝叶斯决策规则进行分类,若 x 位于拒绝域 R ,则拒绝对 x 分类。即决策规则为

$$\text{若 } p(\omega_i | x) = \max_j p(\omega_j | x) \geq 1 - t, \text{ 则 } x \in \omega_i \quad (2-19)$$

$$\text{若 } \max_i p(\omega_i | x) < 1 - t, \text{ 则拒绝对 } x \text{ 分类} \quad (2-20)$$

正确分类概率 $a(t)$ 是阈值 t 的函数,由前面的公式(2-12)给出,只是此时公式中的积分域应改为接受域 A :

$$a(t) = \int_A \max_i [p(x | \omega_i) p(\omega_i)] dx \quad (2-21)$$

拒绝 x 的非条件概率 $r(t)$ 为

$$r(t) = \int_R p(x) dx \quad (2-22)$$

于是,在公式(2-13)中,将积分区域改为A,就得到接受x但对其进行了错误分类的概率为

$$\begin{aligned} e(t) &= \int_A [1 - \max_i p(\omega_i | x)] p(x) dx \\ &= \int_A p(x) dx - \int_A \max_i [p(x | \omega_i) p(\omega_i)] dx = 1 - r(t) - a(t) \end{aligned} \quad (2-23)$$

可以看出,错误率 $e(t)$ 与拒绝率 $r(t)$ 呈负相关的关系。

2.2 最小风险贝叶斯决策规则

上节例2.1中,根据 $p(\text{阳} | \omega_1)p(\omega_1) > p(\text{阳} | \omega_2)p(\omega_2)$ 或 $p(\omega_1 | \text{阳}) > p(\omega_2 | \text{阳})$ 将某人判属正常人 ($\in \omega_1$),并非他属正常人的概率为100%。

实际上,

$$\begin{aligned} p(\omega_1 | \text{阳}) &= \frac{p(\text{阳} | \omega_1)p(\omega_1)}{p(\text{阳} | \omega_1)p(\omega_1) + p(\text{阳} | \omega_2)p(\omega_2)} \\ &= \frac{0.01 \times 0.995}{0.01 \times 0.995 + 0.95 \times 0.005} \approx 67.7\% \end{aligned}$$

$$p(\omega_2 | \text{阳}) = 32.2\%$$

即他属于正常人的概率为67.7%,属于癌症病人的概率为32.2%。因此,仍有可能错判,而错判会造成损失,本为癌症病人错判正常人损失比本是正常人错判为癌症病人的损失大,因为前者将使癌症病人延误治疗,造成生命危险。

现考虑两类问题,引进损失矩阵:

$$(\lambda_{ij})_{2 \times 2} = \begin{bmatrix} \lambda_{11} & \lambda_{12} \\ \lambda_{21} & \lambda_{22} \end{bmatrix} \quad (2-24)$$

其中 $\lambda_{ij} (i, j = 1, 2)$ 为 $x \in \omega_i$ 而判属 ω_j 所造成的损失,称为损失函数。

条件风险定义为将 x 判属某类所造成的损失的条件数学期望,即

模式 x 判属 ω_1 类的条件风险为:

$$r_1(x) = \lambda_{11}p(\omega_1 | x) + \lambda_{21}p(\omega_2 | x) \quad (2-25)$$

模式 x 判属 ω_2 类的条件风险为:

$$r_2(x) = \lambda_{12}p(\omega_1 | x) + \lambda_{22}p(\omega_2 | x) \quad (2-26)$$

最小风险贝叶斯决策规则为

(1) 后验概率形式: