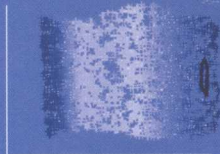
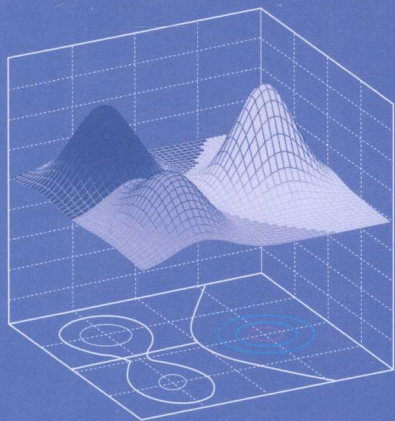


模式识别

Pattern Recognition

孙即祥 姚伟 滕书华 编著



国防工业出版社

National Defense Industry Press

模式识别

孙即祥 姚伟 滕书华 编著

国防工业出版社

·北京·

图书在版编目(CIP)数据

模式识别/孙即祥,姚伟,滕书华编著. —北京:国防工业出版社,2009.11

ISBN 978-7-118-06433-9

I. 模... II. ①孙...②姚...③滕... III. 模式识别
IV. 0235

中国版本图书馆 CIP 数据核字(2009)第 115821 号

※

国防工业出版社出版发行

(北京市海淀区紫竹院南路 23 号 邮政编码 100048)

北京四季青印刷厂印刷

新华书店经售

*

开本 880 × 1230 1/32 印张 9 $\frac{1}{2}$ 字数 260 千字

2009 年 11 月第 1 版第 1 次印刷 印数 1—5000 册 定价 22.00 元

(本书如有印装错误,我社负责调换)

国防书店:(010)68428422

发行邮购:(010)68414474

发行传真:(010)68411535

发行业务:(010)68472764

前 言

分类识别是人们最重要的基本活动之一,在人们的日常生活、社会活动、科研生产以及学习、工作中无时无刻不在进行着分类识别。模式识别是研究分类识别理论和方法的科学,是一门综合性、交叉性学科。在理论上它涉及代数学、矩阵论、函数论、随机数学、模糊数学、图论、最优化理论、信号处理、计算机科学、神经物理学等众多学科的知识,在应用上又与其他许多领域知识及工程技术密切相关。自 20 世纪 80 年代以来,它受到了学术界和各应用领域的极大重视,计算机软、硬件技术的日臻成熟及其他相关学科的迅速发展更使它成为理论研究和技术开发的热们学科。

本书是一本关于模式识别理论和方法的教材,在撰写过程中遵循以下原则:在结构安排上尽量使知识表达体系与学科本身的体系相一致;在内容阐述方式上遵循人的认知规律;在选材上尽量满足读者掌握基础的学科知识。我们的目标就是使本书可读性好、学术性强、实用价值大。本书具有下述的特点。

(1)内容广泛。模式识别是一门相当活跃的重要学科,其发展非常迅速,所涉及的理论十分广泛、方法十分丰富,新理论、新方法、新技术、新应用不断涌现。本书系统地阐述了模式识别领域的基础知识及经典方法,对经实践证明有重要现实意义的新理论、新方法和新技术也进行了介绍。本书包括统计模式识别、模糊模式识别、神经网络技术、人工智能等方法。

(2)结构清晰合理。学科知识的组织结构是否得当将会直接影响读者的学习效果,合理清晰的学科知识表述体系有益于读者对各种理论、方法的理解和记忆。本书处理好模式识别知识点的布局,全书及各章节都尽量按由浅入深、先易后难、先具体后抽象来安排。

(3)选材考究精细。如前所述,这门学科的各种理论、方法、技术纷繁众多,而且还在不断地出现和发展。本书在众多的知识中选取了那些具有基础理论性、思维训练性及有效实用性的重要内容,并且注意处理好本书内容与其他学科知识的关系。

(4)注重基础。打好基础是教育经验的总结,也是科技高速发展的需要,本书自始至终都非常注重强化基本概念、基本思想、基础理论、基本方法和基本技能。

(5)详略得当。由于本书涉及的知识面广,所以必须在论述上有详有略,详略得当。对于统计模式识别等基础的经典内容,详细叙述其基本的、重要的理论和方法;而对于模糊模式识别、神经网络等技术,则重点介绍其思想概念,让读者对其有所了解。这样既实现了知识在面上的宽广,又达到了知识点处的理论深度。

与模式识别直接密切相关的课程是矩阵论、概率论与数理统计、最优化理论与方法等。读者只要具备一些必要的理论基础和相关基本知识便可以顺利地学完每一章的主要内容。对于希望深入学习的读者,则应先行修完上述课程,并根据本书提供的参考文献参阅其他相关的资料。

由于模式识别是一门不断发展的学科,新的理论、方法和技术、新的应用成果不断涌现,再加上我们的学识水平及时间有限,本书可能没有完全达到我们所希望的目标,也不可避免地存在各种错误和疏漏,敬请读者给予批评指正。

作者

2009年7月于国防科技大学

目 录

第 1 章 引论	1
1.1 概述	1
1.1.1 模式识别的概念	1
1.1.2 模式识别系统	2
1.1.3 模式识别的基本方法	4
1.2 特征矢量和特征空间	6
1.3 随机矢量的描述	7
1.4 正态分布	12
1.4.1 正态分布的定义	12
1.4.2 多元正态分布的性质	14
第 2 章 聚类分析及最近邻方法	19
2.1 聚类分析的概念	19
2.1.1 聚类分析的基本思想	19
2.1.2 特征量	20
2.1.3 方法的有效性	20
2.2 模式相似性测度	22
2.2.1 距离测度(差值测度)	23
2.2.2 相似测度	26
2.2.3 匹配测度	28
2.3 类间距离	30

2.3.1	类间距离测度方法	30
2.4	准则函数	32
2.4.1	点与集间的距离	32
2.4.2	聚类的准则函数	35
2.5	聚类的算法	39
2.5.1	聚类的技术方案	39
2.5.2	基于相似性阈值的简单聚类方法	39
2.5.3	谱系聚类法	42
2.5.4	动态聚类法(Dynamic clustering algorithm)	46
2.5.5	近邻函数法	55
2.6	最近邻方法	58
2.6.1	最近邻法	59
2.6.2	剪辑最近邻法	59
2.6.3	引入拒绝类别决策的最近邻法	62
	习题	63
	算法编程	66
第3章 判别域代数界面方程法		69
3.1	判别域界面方程分类的概念	69
3.2	线性判别函数	70
3.2.1	两类问题	71
3.2.2	多类问题	71
3.3	判别函数值的鉴别意义、权空间及解空间	77
3.3.1	判别函数值的大小、正负的数学意义	77
3.3.2	权空间、解矢量与解空间	79
3.4	Fisher 线性判别	81
3.5	线性可分条件下判别函数权矢量算法	86
3.5.1	感知器算法	86

3.5.2	一次准则函数及梯度下降法	89
3.5.3	感知器训练算法在多类问题中的应用	92
3.6	一般情况下的判别函数权矢量算法	94
3.6.1	分段二次准则函数及共轭梯度法	95
3.6.2	最小平方误差准则及 W-H 算法	97
3.6.3	H-K(Ho-Kashyap)算法	99
3.7	广义线性判别函数	102
3.8	二次判别函数	105
3.9	位势函数分类法	106
3.9.1	位势函数的概念	106
3.9.2	由位势函数产生判别函数的训练算法及分类规则 ..	108
3.10	支持矢量机简介	111
	习题	115
	算法编程	116
第4章	统计判决	118
4.1	最小误判概率准则判决	119
4.1.1	最小误判概率准则判决的一般形式	119
4.1.2	正态模式最小误判概率判决规则的具体形式	125
4.1.3	正态模式分类的误判概率	130
4.2	最小损失准则判决	134
4.2.1	损失概念、损失函数与平均损失	134
4.2.2	最小损失准则判决	136
4.3	最小最大损失准则	140
4.4	N-P(Neyman-Pearson)判决	144
4.5	序贯判决(SPRD)	149
4.5.1	控制误判概率的序贯判决	150
4.5.2	计入提取特征代价的最小损失准则下的序贯	

判决	154
习题	155
算法编程	161
第 5 章 统计决策中的经典学习方法	164
5.1 统计推断概述	164
5.2 参数估计	166
5.2.1 均值向量和协方差阵的矩法估计	166
5.2.2 最大似然估计(MLE)	169
5.2.3 贝叶斯估计(BE)	171
5.2.4 最大似然估计和贝叶斯估计的性能比较	173
5.3 贝叶斯学习	174
5.4 概密的窗函数估计法	178
5.4.1 概密的基本估计式	178
5.4.2 Parzen 窗法	180
5.4.3 k_N -近邻估计法	181
5.4.4 后验概率的估计	182
5.5 有限项正交函数级数逼近法	183
5.5.1 最小积分平方差逼近方法	183
5.5.2 最小均方差逼近方法	186
5.6 用位势函数法逼近贝叶斯判决函数	187
5.7 错误率估计	189
5.7.1 分类器错误率的实验估算基本原理	189
5.7.2 样本抽取方式对误判概率估计的影响	190
5.7.3 训练与测试样本集的大小对错误率的影响	191
5.7.4 训练样本使用技术及错误率的测试方法	192
习题	193
算法编程	197

第 6 章 特征提取与选择	201
6.1 概述	201
6.2 类别可分性判据	202
6.2.1 基于几何距离的可分性判据	203
6.2.2 基于类的概率密度函数的可分性判据	207
6.2.3 基于后验概率的可分性判据	209
6.3 基于可分性判据进行变换的特征提取	211
6.3.1 基于离差阵的特征提取	211
6.3.2 多类情况	213
6.3.3 基于熵概念的某些特征提取与选择方法	214
6.4 最佳鉴别矢量的提取	214
6.4.1 Fisher 鉴别矢量及鉴别平面	215
6.4.2 最佳鉴别矢量集	217
6.5 离散 K-L 变换及其在特征提取与选择中的应用	218
6.5.1 离散 K-L 变换(DKLT)	219
6.5.2 离散 K-L 变换在特征提取与选择中的应用	223
6.6 特征选择中的直接挑选法	230
6.6.1 次优搜索法	230
6.6.2 最优搜索法	233
习题	235
第 7 章 其他模式识别方法	238
7.1 模糊模式识别	238
7.1.1 模糊数学基础知识	238
7.1.2 模糊模式识别的基本方法	244
7.2 神经网络在模式识别中的应用	246
7.2.1 人工神经网络的基本知识	246

7.2.2 常见的神经网络模型	250
7.3 句法模式识别	253
7.3.1 句法模式识别概述	253
7.3.2 形式语言介绍	254
7.3.3 句法模式识别的基本内容	259
7.4 信息融合	260
7.4.1 信息融合概述	260
7.4.2 融合技术层次性及融合系统功能模块和结构	261
7.5 树分类器	268
7.5.1 树分类器原理及设计原则	269
7.5.2 树分类器关键技术	270
7.5.3 决策树生成算法	271
习题	276
参考文献	278

第 1 章 引 论

1.1 概 述

1.1.1 模式识别的概念

人们在日常生活、社会活动、科研生产以及学习、工作中无时无刻不在进行着分类识别,分类识别是人类的基本活动之一。例如,儿童在认读识字卡片上的数字时,将它们区分为 0~9 中的一个,这是对数字符号的识别;在读书看报时,人们进行的是文字识别活动;做某种实验时对示波器显示波形的观察是一个波形识别过程;医生给患者诊断疾病需要对病情进行识别;在人群中寻找某一个人是对人的形体及其他特征的识别行为。随着人类社会活动及生产科研广泛而深入的发展,需要识别的对象种类越来越多,内容越来越深入和复杂,要求也越来越高,为了改善工作条件、减轻工作强度,人们希望机器能代替人类完成某些繁重的识别工作;有些场合环境恶劣、存在危险或人们根本就不能接近,这就需要借助机器、运用分析算法进行识别。人们利用机器可以提高识别的速度、正确率及扩大应用的广度。我们这里所说的模式识别是指运用机器进行分类识别。模式识别的重要应用之一是计算机自动诊断疾病。与医生的诊断过程相仿,首先要获得就诊人的有关情况,如要测量体温、血压、心率,还可能对血液等进行化验,作 X 光透视、B 超、心电图,甚至 CT 等。医生根据这些量测结果以及病人的病史等资料,运用自己的临床经验对患者进行诊断。而机器识别是将上述各种有用的资料信息输入计算机中,在这之前计算机已装入有关的分析算法,这些算法是专家知识、经验的总结和集成,其形式可以是规则、函数、数表等,通过计算机程序对信息进行分析并作出判断。

下面首先较详细地介绍模式和模式识别这两个基本概念。为了能

让机器执行和完成分类识别任务,必须首先将分类识别对象的有用的信息输入计算机中,为此,应对分类识别对象进行科学的抽象,建立它的数学模型,用以描述和代替识别对象,称这种对象的描述为模式 (Pattern)。Pattern 的原义是模范、模型、典型、样品、图案等,其内涵深刻、外延广泛。无论是自然界中物理、化学或生物等领域的对象,还是社会中的语言、文字等,都可以进行科学的抽象,具体地讲,我们对它们进行量测,得到表征它们特征的一组数据,为使用方便,将它们表示成矢量形式,称其为特征矢量;也可以将对象的特征属性作为基元、用符号表示,从而将它们的结构特征描述成一个符号串、图或某个数学式子。通俗地讲,模式就是事物的代表,是事物的数学模型之一,它的表示形式是矢量、符号串、图或数学关系。对一类对象的抽象也称为该类的模式。

所谓模式识别 (Pattern Recognition),是根据研究对象的特征或属性,利用以计算机为中心的机器系统,运用一定的分析算法认定它的类别,系统应使分类识别的结果尽可能地符合真实。

目前,模式识别理论和技术已成功地应用于工业、农业、金融、军事、公安、科研、生物医学、气象、天文学等许多领域,如我们熟知的信件分拣、遥感图片的机器判读、系统的故障诊断、生物特征识别(指纹、虹膜、脸等识别)、生物医学的细胞或组织分析、具有视觉的机器人、武器制导寻的系统、汽车自动驾驶系统以及文字与语言的识别等,并且现在正扩展到许多其他领域。当今时代科技发展的重要趋势之一是智能化,模式识别是人工智能的一个分支。尽管现在机器识别的水平还远不如人脑,但随着模式识别理论以及其他相关学科的发展,可以预言,它的功能将会越来越强大,应用也会越来越广泛。

1.1.2 模式识别系统

一个功能较完善的识别系统在进行模式识别之前,首先需要进行学习。模式识别系统及识别过程的原理框图可以用图 1.1.1 表示。虚线的上部是分类、识别过程,虚线的下部是学习、训练过程。需要指出的是,应用的目的不同、采用的分类识别方法不同,具体的分类识别系统和过程将会有所不同。一般而言,特征提取与选择、训练学习、分类

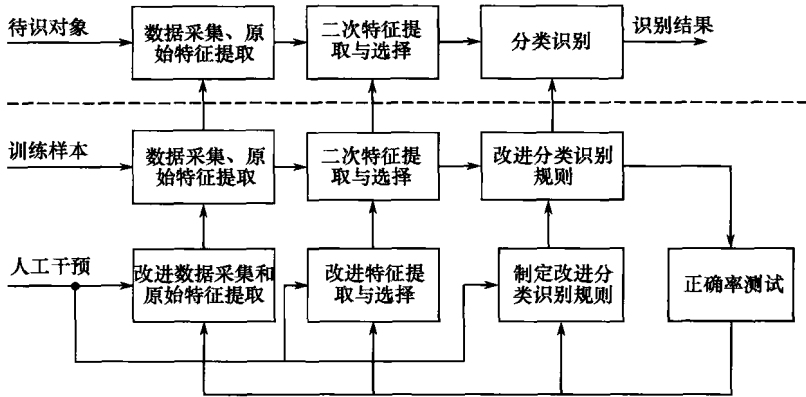


图 1.1.1 模式识别系统及识别过程的原理框图

识别是任何模式识别方法或系统的三大核心问题。

下面对识别系统的主要环节作简要的说明。

1. 特征提取

无论是识别过程还是学习过程,都要对研究对象固有的、本质的及重要的特征或属性进行量测并将结果数值(字)化,或将对象分解并符号化,形成特征矢量或符号串、关系图,从而产生代表对象的模式,模式类中的个体在有些场合中也称为样本。用于学习与训练的样本的类别通常是已知的。另外,在进行特征提取之前一般还需要对目标的有关信息进行预处理。

2. 特征选择

通常能描述对象的特征或属性的种类很多,为了节约资源,节省计算机存储空间、运算时间、特征提取的费用,有时更是为了算法的可行性,在满足分类识别正确率要求的条件下,按某种准则尽量选用对正确分类识别作用较大的特征,使得用较少的特征就能完成分类识别任务。这项工作表现为减少特征矢量的维数、符号串字符数或简化图的结构。

3. 学习和训练

为了让机器具有分类识别功能,如同人类自身一样,人们应首先对它进行训练,将人类的识别知识和方法以及关于分类识别对象的知识输入机器中,产生分类识别的规则和分析程序,这也相当于机器进行学

习。这个过程一般要反复进行多次,不断地修正错误、改进不足,这包括修正特征提取方法、特征选择方案、判决规则方法及参数,最后使系统正确识别率达到设计要求。目前,机器学习常需要人工干预,这个过程通常是人机交互的。

4. 分类识别

在学习、训练之后,所产生的分析规则及程序用于未知类别的对象的分类识别。需要指出的是,输入机器的人类分类识别的知识和方法以及有关对象知识越充分,这个系统的识别功能就越强、正确率就越高,有些分类过程(如聚类分析)似乎没有将有关对象的知识输入,实际上我们在选择距离测度、在采用某种聚类方法时,已经隐含地用到了对象的一些知识,也在一定程度上加入了人们的知识。

1.1.3 模式识别的基本方法

由于分类识别活动是人们的基本活动,人们希望机器能代替人们进行分类识别工作,因此模式识别的理论和方法引起人们极大的兴趣并进行了长期的深入研究,现已发展成一门多学科的交叉学科。这门学科涉及的理论与技术相当广泛,包括多种数学理论、神经心理学、计算机科学、信号处理等。从本质上讲,这门学科实际上是数据处理与信息分析;从功能上讲,可以认为它是人工智能的一个分支。

针对不同的对象和不同的目的,可以运用不同的模式识别理论、方法,目前主流的技术是统计模式识别、句法模式识别、模糊模式识别、人工神经网络方法、人工智能方法、子空间方法,它们之间往往存在一定的联系和借鉴。

1. 统计模式识别

这类识别技术理论较完善,方法也很多,通常较为有效,现已形成了一个完整的体系。尽管方法很多,但从根本上讲,都是直接利用各类的分布特征,即利用各类的概率密度函数、后验概率等,或隐含地利用上述概念进行分类识别。其中基本的技术为聚类分析法、判别类域代数界面法、统计决策法、最近邻法等。在聚类分析中,利用待分类模式之间的“相似性”进行分类,较相似的作为一类,较不相似的不作为一类。在分类过程中不断地计算所分划的各类的中心,待分类模式与各

类中心的距离作为对其分类的依据。这实际上是在某些设定下隐含地利用了概率分布概念,因常见的概率密度函数中,距期望值较近的点概率密度较大。该类方法的另一种技术是根据待分类模式和已指判出类别的模式之间的距离来确定其类别,这实际上也是在一定程度上利用了有关的概念。最近邻法是根据待分类模式的一个或 k 个近邻训练样本的类别而确定其类别。判别类域界面法中,用已知类别的训练样本产生判别函数,这相当于学习或训练,根据待分类模式代入判别函数后所得值的正负确定其类别,判别函数提供了相邻两类判别域的界面,其也相应于在一些设定下两类概率密度函数之差。在统计判决中,在一些分类识别准则下严格地按照概率统计理论导出各种判决规则,这些判决规则可以产生某种意义上的最优分类识别结果,这些判决规则要用到各类的概率密度函数、先验概率或后验概率。这可以通过训练样本对未知概率密度函数中的参数进行估计,或对未知的概率密度函数进行逼近而估计它们。

2. 句法模式识别

句法模式识别也称为结构模式识别。在许多情况下,对于较复杂的对象仅用一些数值特征已不能较充分地对其描述与正确识别,这时可采用句法识别技术。句法识别技术是将对象分解成若干个基本单元,这些基本单元称为基元,用这些基元以及它们的结构关系来表征对象,基元以及这些基元的结构关系可以用字符串或图来表示,这些字符串或图进一步抽象为语言的句子,然后根据代表类的文法运用形式语言理论对该句子进行句法分析,根据其是否符合某一类的文法而决定其类别。

3. 模糊模式识别

这类识别技术运用模糊数学的理论和解决方法解决模式识别问题,因此适用于分类识别对象本身或要求的识别结果具有模糊性的场合。模糊模式识别方法的基本思想是,将模式或模式类作为模糊集,将模式的属性值或属性转化为隶属度,运用隶属度、模糊关系或模糊推理进行分类识别。目前,模糊识别方法较多。这类方法的有效性主要在于对象类的隶属函数建立得是否良好,对象间的模糊关系的度量是否良好。

4. 人工神经网络方法

人工神经网络是由大量简单的基本单元——神经元(Neuron)相

互联结而构成的非线性动态系统,每个神经元结构和功能比较简单,而由其构成的系统却可以非常复杂,具有生物神经网络的某些特性,在自学习、自组织、联想及容错方面具有较强的能力,可用于联想、识别和决策。在模式识别方面,与前述方法显著不同的特点之一是在学习过程中具有自动提取特征的能力。

5. 人工智能方法

众所周知,人类具有极完善的分类识别功能,人工智能是研究如何使机器具有人类智能的理论和方法,模式识别从本质上讲就是如何根据对象的特征进行类别的判断,因此可将人工智能中有关知识表示、推理、学习等技术用于模式识别。

6. 子空间方法

子空间方法是将代数学的基本理论与统计基本理论进行综合,这类方法的基本思想是,根据各类训练样本的相关阵通过线性变换由原始模式特征空间产生各类对应的子空间,这些子空间由原始模式类的样本相关阵的主要特征矢量所张成,这些主要特征矢量反映了模式分布结构信息,每个子空间与每个类别一一对应。在子空间方法中,主要的分类决策规则有基于投影长度的比较方法、基于表达熵的比较方法和基于统计假设检验的方法三个类型。基于投影的识别方法是根据待识模式在各个子空间的投影大小判定该模式类别。

上述的6类方法各有特点及其应用范围,现在来看,它们不能相互取代,只能共存,相互促进、借鉴、渗透及融合。除了上述的6类方法外,还有其他的一些类型方法,如协同模式识别方法等。一个较完善的识别系统很可能是综合利用上述各类识别方法的观点、概念和技术而构成的。

1.2 特征矢量和特征空间

设一个研究对象的 n 个特征量测值分别为 x_1, x_2, \dots, x_n ,由于它们来自同一个对象,所以应将它们作为一个整体一起考虑,为此,让它们构成一个 n 维特征矢量 \mathbf{x} ,即 $\mathbf{x} = (x_1, x_2, \dots, x_n)^T$, \mathbf{x} 是原对象的一种数学抽象,用其来代表原对象,即为原对象的模式。对某对象的分类识别