

模式识别与 智能计算

——Matlab技术实现

● 杨淑莹 著



含光盘1张



电子工业出版社
PUBLISHING HOUSE OF ELECTRONICS INDUSTRY

<http://www.phei.com.cn>

模式识别与智能计算 ——Matlab 技术实现

杨淑莹 著

電子工業出版社

Publishing House of Electronics Industry

北京·BEIJING

内 容 简 介

本书广泛吸取统计学、神经网络、数据挖掘、机器学习、人工智能、群智能计算等学科的先进思想和理论,将其应用到模式识别领域中;以一种新的体系,系统、全面地介绍模式识别的理论、方法及应用。全书共分为13章,内容包括:模式识别概述,特征的选择与提取,模式相似性测度,贝叶斯分类器设计,判别函数分类器设计,神经网络分类器设计(BP神经网络、径向基函数神经网络、自组织竞争神经网络、概率神经网络、对向传播神经网络、反馈型神经网络),决策树分类器,粗糙集分类器,聚类分析,模糊聚类分析,遗传算法聚类分析,蚁群算法聚类分析,粒子群算法聚类分析。

本书内容新颖,实用性强,理论与实际应用密切结合,以手写数字识别为应用实例,介绍理论运用于实践的实现步骤及相应的Matlab代码,为广大研究工作者和工程技术人员对相关理论的应用提供借鉴。

本书可作为高等院校计算机工程、信息工程、生物医学工程、智能机器人学、工业自动化、模式识别等学科本科生、研究生的教材或教学参考书,亦可供相关工程技术人员参考。

未经许可,不得以任何方式复制或抄袭本书之部分或全部内容。
版权所有,侵权必究。

图书在版编目(CIP)数据

模式识别与智能计算: Matlab 技术实现/杨淑莹著. —北京:电子工业出版社, 2008.1
ISBN 978-7-121-05453-2

I. 模… II. 杨… III. ①模式识别-计算机辅助计算-软件包, MATLAB ②人工智能-计算机辅助计算-软件包, MATLAB IV. O235-39 TP183

中国版本图书馆 CIP 数据核字(2007)第 181139 号

责任编辑:张 榕 文字编辑:毕军志

印 刷:北京市海淀区四季青印刷厂

装 订:涿州市桃园装订有限公司

出版发行:电子工业出版社

北京市海淀区万寿路 173 信箱 邮编 100036

开 本:787×1092 1/16 印张:22.75 字数:580 千字

印 次:2008 年 1 月第 1 次印刷

印 数:5000 册 定价:48.00 元(含光盘 1 张)

凡所购买电子工业出版社的图书,如有缺损问题,请向购买书店调换。若书店售缺,请与本社发行部联系,联系及邮购电话:(010) 88254888。

质量投诉请发邮件至 zltz@phei.com.cn, 盗版侵权举报请发邮件至 dbqq@phei.com.cn。

服务热线:(010) 88258888。

反侵权盗版声明

电子工业出版社依法对本作品享有专有出版权。任何未经权利人书面许可,复制、销售或通过信息网络传播本作品的行为;歪曲、篡改、剽窃本作品的行为,均违反《中华人民共和国著作权法》,其行为人应承担相应的民事责任和行政责任,构成犯罪的,将被依法追究刑事责任。

为了维护市场秩序,保护权利人的合法权益,我社将依法查处和打击侵权盗版的单位和个人。欢迎社会各界人士积极举报侵权盗版行为,本社将奖励举报有功人员,并保证举报人的信息不被泄露。

举报电话:(010)88254396;(010)88258888

传 真:(010)88254397

E-mail: dbqq@phei.com.cn

通信地址:北京市万寿路173信箱

电子工业出版社总编办公室

邮 编:100036

读者调查表

尊敬的读者：

欢迎您参加读者调查活动,对我们的图书提出真诚的意见,您的建议将是我们创造精品的动力源泉。为方便大家,我们提供了两种填写调查表的方式:

1. 您可以登录 <http://yydz.phei.com.cn>,进入右上角的读书栏目,填好本调查表后直接反馈给我们。
2. 您可以填写下表后寄给我们(北京海淀区万寿路 173 信箱应用电子技术图书事业部 邮编:100036)。

姓名: _____ 性别: 男 女 年龄: _____ 职业: _____
电话(寻呼): _____ E-mail: _____
传真: _____ 通信地址: _____
邮编: _____

1. 影响您购买本书的因素(可多选):

- 封面封底 价格 内容简介、前言和目录 书评广告 出版物名声
作者名声 正文内容 其他 _____

2. 您对本书的满意度:

- 从技术角度 很满意 比较满意 一般 较不满意 不满意
从文字角度 很满意 比较满意 一般 较不满意 不满意
从排版、封面设计角度 很满意 比较满意 一般 较不满意
不满意

3. 您最喜欢书中的哪篇(或章、节)? 请说明理由。

4. 您最不喜欢书中的哪篇(或章、节)? 请说明理由。

5. 您希望本书在哪些方面进行改进?

6. 您感兴趣或希望增加的图书选题有:

邮寄地址:北京海淀区万寿路 173 信箱应用电子技术分社 张榕 收 邮编:100036

编辑电话:(010)88254455 E-mail:zr@phei.com.cn

前 言

模式识别成为当代高科技研究的重要领域之一,发展成为一门独立的新学科。模式识别技术迅速扩展,已经应用在人工智能、机器人、系统控制、遥感数据分析、生物医学工程、军事目标识别等领域,几乎遍及各个学科领域,在国民经济、国防建设、社会发展的各个方面得到广泛应用,产生深远的影响。

本书的宗旨

国内外论述模式识别技术的参考书为数不少,这一领域涉及深奥的数学理论,往往使实际工作者感到困难,大部分书只是罗列模式识别的各种算法,看不到算法的实际效果和各种算法对比的结果,而这正是学习者和实际工作者所需要了解和掌握的内容。目前还确实缺少一本关于模式识别技术在实际应用方面具有系统性、可比性和实用性的参考书。

本书内容基本涵盖了目前模式识别重要的理论和方法,但并没有简单地将各种理论方法堆砌起来,而是将作者自身的研究成果和实践经验传授给读者,在介绍各种理论和方法的同时,将不同算法应用于实际中,书中含有需要应用模式识别技术解决的问题、模式识别理论的讲解和推理、将理论转化为编程的步骤、计算机能够运行的源代码、计算机运行模式识别算法程序后的效果、不同算法应用于同一个问题的效果对比;使读者不至于面对如此丰富的理论和方法无所适从,而是有所学即有所用。

由于至今还没有统一的、有效的能够应用于所有的模式识别的理论,当前的一种普遍看法认为,不存在对所有的模式识别问题都适用的单一模型和解决识别问题的单一技术,我们所要做的是把模式识别方法与具体问题结合起来,把模式识别与统计学、神经网络、数据挖掘、机器学习、人工智能、群智能计算等学科的先进思想和理论结合起来,为读者提供一个多种理论的测试平台,并在此基础上,深入掌握各种理论的效能和应用的可能性,互相取长补短,开创模式识别应用的新局面。

本书的主要内容

本书共 13 章,分成四部分,主要介绍统计模式识别。由于篇幅有限,没有讨论句法模式识别。

第 1 章介绍模式识别的基本概念。

第 2 章介绍特征的选择与提取。

第 3 章至第 8 章介绍分类器设计问题。

➤ 第 3 章介绍模式相似性测度。

➤ 第 4 章介绍基于概率统计的贝叶斯分类器设计。

➤ 第 5 章介绍判别函数分类器设计。

➤ 第 6 章介绍神经网络分类器设计。

- 第7章介绍决策树分类器设计。
- 第8章介绍粗糙集分类器设计。
- 第9章至第13章介绍聚类分析问题。
- 第9章介绍聚类分析。
- 第10章介绍模糊聚类分析。
- 第11章介绍遗传算法聚类分析。
- 第12章介绍蚁群算法聚类分析。
- 第13章介绍粒子群算法聚类分析。

本书的特点

1. 选用新技术。除了介绍许多重要经典的内容以外,书中还包括了最近十几年来刚刚发展起来的并被实践证明有用的新技术、新理论,例如,支持向量机、BP神经网络、径向基函数神经网络、自组织竞争神经网络、概率神经网络、对向传播神经网络、反馈型神经网络、决策树、粗糙集理论、模糊集理论、模拟退火、遗传算法、蚁群算法、粒子群算法等,并将这些新技术应用于模式识别当中,提供这些新技术的实现方法和源代码。

2. 实用性强,针对实例介绍理论和技术,使理论和实践相结合,避免了空洞的理论说教。书中实例取材于手写数字模式识别,对于数字识别属于多类问题,在实际应用中具有广泛的代表性,读者对程序稍加改进,就可以应用到不同的场合,例如,文字识别、字符识别、图形识别等。

3. 针对每一种模式识别技术,书中分为理论基础、实现步骤、编程代码三部分;在掌握了基本理论之后,按照实现步骤的指导,可以了解算法的实现思路和方法,再进一步体会短小精悍的核心代码,学习者可以很快掌握模式识别技术,经过应用本书提供的实例程序,立刻会见到算法的实际效果。所有算法都用 Matlab 编程实现,便于读者学习和应用。

本书的读者对象

本书可作为高等院校计算机工程、信息工程、生物医学工程、智能机器人学、工业自动化、模式识别等学科研究生、本科生的教材或教学参考书,亦可供有关工程技术人员参考。

本书的光盘

本书配套光盘中提供了两个综合性的实例软件:分类识别软件和聚类分析软件,每个软件包含了多种模式识别算法。

1. 分类识别软件。分类识别软件界面、菜单调用功能和运行效果如图 0-1 所示。

当读者运行本书配套光盘分类识别软件时,首先随机手写 0~9 的数字,手写数字的背景是白色,笔是黑色,当用户不满意时可以删除该数字。然后单击菜单,调用各种模式识别算法,输出识别结果。针对同一次手写的数字,应用各种算法进行识别,读者可以比较各种算法的识别效果。

2. 聚类判别分析软件。聚类判别分析软件界面、菜单调用功能和运行效果如图 0-2 所示。

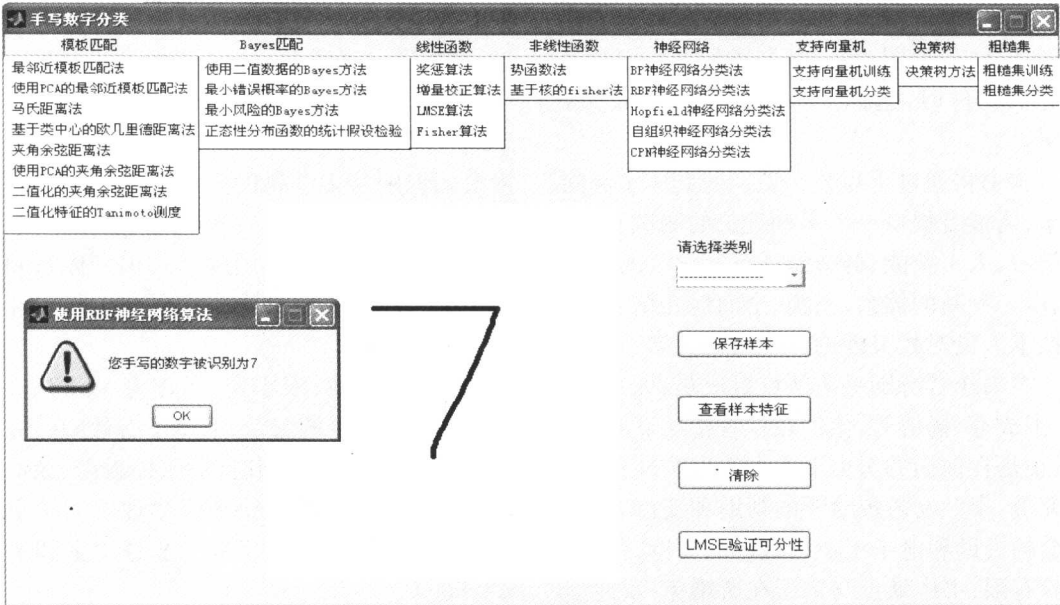


图 0-1 分类识别软件界面、菜单调用功能和运行效果图



图 0-2 聚类判别分析软件界面、菜单调用功能和运行效果图

聚类判别分析软件有两种输入数据方式:一种方式是打开一幅含有需要聚类分析的 8 位位图或 24 位位图;另一种方式是直接用画笔画出各种图形。当读者运行本书配套光盘聚类判别分析软件时,选择一种输入数据方式,然后单击菜单,调用各种模式识别算法,输出聚类结果。

本书将通过手写数字识别的具体实例问题,将模式识别与多学科的先进思想和理论结合起来,为读者提供一个多种理论的测试平台。本书广泛吸取了统计学、神经网络、数据挖掘、机器学习、人工智能、群智能计算等学科的先进思想和理论,将其扩充应用到模式识别体系领域中,以一种新的体系,系统、全面地介绍模式识别的理论、方法及应用,为广大研究工作者和工程技术人员对相关理论的应用提供借鉴。

本书作者特别感谢项目组成员:张成、王立群、安博、任翠池、宋小寅、杨作寿、邓和平、杨松、王成芬、杨倩等,他们在作者指导下的研究工作中付出了辛苦的劳动,取得了有效的研究成果,正是在他们的努力下本书得以顺利完成,在此表示衷心的感谢。同时,对张桦教授、徐伯夏研究员、李兰友教授给予的帮助和支持表示衷心的感谢。本书的出版得到天津理工大学出版基金的资助和电子工业出版社的支持,在此作者表示由衷的感谢。由于编者业务水平和实践经验有限,书中缺点与错误在所难免,欢迎读者予以指正!

作者将不辜负广大读者的期望,努力工作,不断充实新的内容。为方便广大读者、用户使用本书配套提供的软件,及时解答用户在阅读、使用中遇到的问题,提供了技术支持电子邮箱:ysying1262@126.com。读者可通过该邮箱及时与作者取得联系,获得技术支持。

作者 谨识

目 录

第 1 章 模式识别概述	1
1.1 模式识别的基本概念	1
1.2 特征空间优化设计问题	4
1.3 分类器设计	6
1.3.1 分类器设计基本方法	8
1.3.2 判别函数	10
1.3.3 分类器的选择	12
1.3.4 训练与学习	13
1.4 聚类设计	13
1.5 模式识别的应用	15
本章小结	15
习题 1	16
第 2 章 特征的选择与提取	17
2.1 样本特征库初步分析	18
2.2 样品筛选处理	19
2.3 特征筛选处理	19
2.3.1 特征相关分析	19
2.3.2 特征选择及搜索算法	20
2.4 特征评估	26
2.5 基于主成分分析的特征提取	29
2.6 特征空间描述与分析	32
2.6.1 特征空间描述	32
2.6.2 特征空间分布分析	37
2.7 手写数字特征提取与分析	40
2.7.1 手写数字特征提取	40
2.7.2 手写数字特征空间分布分析	41
本章小结	45
习题 2	46
第 3 章 模式相似性测度	47
3.1 模式相似性测度的基本概念	47
3.2 距离测度分类法	50
3.2.1 模板匹配法	50
3.2.2 基于 PCA 的模板匹配法	52

3.2.3	基于类中心的欧式距离法分类	54
3.2.4	马氏距离分类	56
3.2.5	夹角余弦距离分类	58
3.2.6	二值化的夹角余弦距离法分类	59
3.2.7	二值化的 Tanimoto 测度分类	60
	本章小结	62
	习题 3	62
第 4 章	基于概率统计的贝叶斯分类器设计	63
4.1	贝叶斯决策的基本概念	63
4.1.1	贝叶斯决策所讨论的问题	63
4.1.2	贝叶斯公式	64
4.2	基于最小错误率的贝叶斯决策	66
4.3	基于最小风险的贝叶斯决策	69
4.4	贝叶斯决策比较	71
4.5	基于二值数据的贝叶斯分类实现	72
4.6	基于最小错误率的贝叶斯分类实现	75
4.7	基于最小风险的贝叶斯分类实现	78
	本章小结	81
	习题 4	82
第 5 章	判别函数分类器设计	83
5.1	判别函数的基本概念	83
5.2	线性判别函数	84
5.3	线性判别函数的实现	88
5.4	感知器算法	89
5.5	增量校正算法	96
5.6	LMSE 验证可分性	102
5.7	LMSE 分类算法	108
5.8	Fisher 分类	111
5.9	基于核的 Fisher 分类	114
5.10	线性分类器实现分类的局限	121
5.11	非线性判别函数	123
5.12	分段线性判别函数	125
5.13	势函数法	128
5.14	支持向量机	133
	本章小结	139
	习题 5	139
第 6 章	神经网络分类器设计	140
6.1	神经网络的基本原理	140
6.1.1	人工神经元	140

6.1.2	人工神经网络模型	143
6.1.3	神经网络的学习过程	146
6.1.4	人工神经网络在模式识别问题上的优势	146
6.2	BP神经网络	147
6.2.1	BP神经网络的基本概念	147
6.2.2	BP神经网络分类器设计	153
6.3	径向基函数神经网络(RBF)	163
6.3.1	径向基函数神经网络的基本概念	163
6.3.2	径向基函数神经网络分类器设计	168
6.4	自组织竞争神经网络	170
6.4.1	自组织竞争神经网络的基本概念	171
6.4.2	自组织竞争神经网络分类器设计	173
6.5	概率神经网络(PNN)	176
6.5.1	概率神经网络的基本概念	176
6.5.2	概率神经网络分类器设计	176
6.6	对向传播神经网络(CPN)	179
6.6.1	对向传播神经网络的基本概念	179
6.6.2	对向传播神经网络分类器设计	181
6.7	反馈型神经网络(Hopfield)	185
6.7.1	Hopfield网络的基本概念	185
6.7.2	Hopfield神经网络分类器设计	188
	本章小结	190
	习题6	190
第7章	决策树分类器	191
7.1	决策树的基本概念	191
7.2	决策树分类器设计	192
	本章小结	199
	习题7	199
第8章	粗糙集分类器	200
8.1	粗糙集理论的基本概念	200
8.2	粗糙集在模式识别中的应用	205
8.3	粗糙集分类器设计	209
	本章小结	222
	习题8	223
第9章	聚类分析	224
9.1	聚类的设计	224
9.2	基于试探的未知类别聚类算法	227
9.2.1	最临近规则的试探法	228
9.2.2	最大最小距离算法	231

9.3	层次聚类算法	234
9.3.1	最短距离法	235
9.3.2	最长距离法	238
9.3.3	中间距离法	242
9.3.4	重心法	245
9.3.5	类平均距离法	249
9.4	动态聚类算法	253
9.4.1	K 均值算法	253
9.4.2	迭代自组织的数据分析算法(ISODATA)	257
9.5	模拟退火聚类算法	262
9.5.1	模拟退火的基本概念	262
9.5.2	基于模拟退火思想的改进 K 均值聚类算法	265
	本章小结	272
	习题 9	272
第 10 章	模糊聚类分析	273
10.1	模糊集的基本概念	273
10.2	模糊集运算	275
10.2.1	模糊子集运算	275
10.2.2	模糊集运算性质	277
10.3	模糊关系	277
10.4	模糊集在模式识别中的应用	282
10.5	基于模糊的聚类分析	283
	本章小结	297
	习题 10	297
第 11 章	遗传算法聚类分析	298
11.1	遗传算法的基本概念	298
11.2	遗传算法的构成要素	300
11.2.1	染色体的编码	300
11.2.2	适应度函数	301
11.2.3	遗传算子	302
11.3	控制参数的选择	304
11.4	基于遗传算法的聚类分析	305
	本章小结	318
	习题 11	318
第 12 章	蚁群算法聚类分析	319
12.1	蚁群算法的基本概念	319
12.2	聚类数目已知的蚁群聚类算法	322
12.3	聚类数目未知的蚁群聚类算法	331
	本章小结	335

习题 12	336
第 13 章 粒子群算法聚类分析	337
13.1 粒子群算法的基本概念	337
13.2 基于粒子群算法的聚类分析	340
本章小结	345
习题 13	346
参考文献	347

第 1 章 模式识别概述

本章要点:

- ☑ 模式识别的基本概念
- ☑ 特征空间优化设计问题
- ☑ 分类器设计
- ☑ 聚类设计
- ☑ 模式识别的应用

1.1 模式识别的基本概念

模式识别(Pattern Recognition)就是机器识别、计算机识别或机器自动识别,目的在于让机器自动识别事物。例如,手写数字的识别,结果就是将手写的数字分到具体的数字类别中;智能交通管理系统的识别,就是判断是否有汽车闯红灯,闯红灯的汽车车牌号码;还有文字识别、语音识别、图像中物体识别,等等。该学科研究的内容是使机器能做以前只能由人类才能做的事,具备人所具有的对各种事物与现象进行分析、描述与判断的部分能力。模式识别是直观的、无所不在的,实际上人类在日常生活的每个环节,都从事着模式识别的活动。人和动物较容易做到模式识别,但对计算机来说却是非常困难的。让机器能识别、分类,就需要研究识别的方法,这就是这门学科的任务。

模式识别研究的目的是利用计算机对物理对象进行分类,在错误概率最小的条件下,使识别的结果尽量与客观物体相符合。机器辨别事物最基本的方法是计算,原则上讲是对计算机要分析的事物与标准模板的相似程度进行计算。例如,要识别一个手写的数字,就要将它与从0~9的模板做比较,看跟哪个模板最相似,或最接近。因此首先要能从度量中看出不同事物之间的差异,才能分辨当前要识别的事物,因此最关键的是找到有效地度量不同类事物的差异的方法。

在模式识别学科中,就“模式”与“模式类”而言,模式类是一类事物的代表,而“模式”则是某一事物的具体体现,例如,数字0、1、2、3、4、5、6、7、8、9是模式类,而用户任意手写的一个数字或任意一个印刷数字则是“模式”,是数字的具体化。

1. 模式的描述方法

在模式识别技术中,被观测的每个对象称为样品,例如,在手写数字识别中,每个手写数字可以作为一个样品,如果共写了 N 个数字,我们把这 N 个数字叫做 N 个样品($X_1, X_2, \dots, X_j, \dots, X_N$),其中0表示有 N_0 个样品,1表示有 N_1 个样品,2表示有 N_2 个样品,3表示有 N_3 个样品,……,一共有 $\omega_1, \omega_2, \dots, \omega_M (M=10)$ 个不同的类别。

对于一个样品来说,必须确定一些与识别有关的因素,作为研究的根据,每一个因素称为一个特征。模式就是样品所具有特征的描述。模式的特征集又可用处于同一个特征空间的特

征向量表示。特征向量的每个元素称为特征,该向量也因此称为特征向量,一般我们用小写英文字母 x, y, z 来表示特征。如果一个样品 \mathbf{X} 有 n 个特征,则可把 \mathbf{X} 看做一个 n 维列向量,该向量 \mathbf{X} 称为特征向量,记做

$$\mathbf{X} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} = (x_1, x_2, \dots, x_n)^T$$

若有一批样品共有 N 个,每个样品有 n 个特征,这些数值可以构成一个 n 行 N 列的矩阵,称为原始资料矩阵,如表 1-1 所示。

表 1-1 原始资料矩阵

特征 \ 样品	\mathbf{X}_1	\mathbf{X}_2	...	\mathbf{X}_j	...	\mathbf{X}_N
x_1	x_{11}	x_{12}	...	x_{1j}	...	x_{1N}
x_2	x_{21}	x_{22}	...	x_{2j}	...	x_{2N}
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
x_i	x_{i1}	x_{i2}	...	x_{ij}	...	x_{iN}
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
x_n	x_{n1}	x_{n2}	...	x_{nj}	...	x_{nN}

模式识别问题就是根据 \mathbf{X} 的 n 个特征来判别模式 \mathbf{X} 属于 $\omega_1, \omega_2, \dots, \omega_M$ 类中的哪一类。待识别的不同模式都在同一特征空间中考察,不同模式类由于性质上的不同,它们在各特征取值范围内有所不同,因而会在特征空间的不同区域中出现。要记住向量的运算是建立在各个分量基础之上的。

因此,模式识别系统的目标是在特征空间和解释空间之间找到一种映射关系。特征空间由从模式得到的对分类有用的度量、属性或基元构成的空间。解释空间由 M 个所属类别的集合构成。

如果一个对象的特征观察值为 $\{x_1, x_2, \dots, x_n\}$,它可构成一个 n 维的特征向量值 \mathbf{X} ,即

$$\mathbf{X} = (x_1, x_2, \dots, x_n)^T$$

式中, x_1, x_2, \dots, x_n 为特征向量 \mathbf{X} 的各个分量。

一个模式可以看做 n 维空间中的向量或点,此空间称为模式的特征空间 R^n 。在模式识别过程中,要对许多具体对象进行测量,以获得许多观测值,其中有均值、方差、协方差与协方差矩阵等。

2. 模式识别系统

一个典型的模式识别系统如图 1-1 所示,由数据获取、预处理、特征提取、分类决策及分类器设计五部分组成。一般分为上下两部分:上半部分完成未知类别模式的分类;下半部分属于分类器设计的训练过程,利用样品进行训练,确定分类器的具体参数,完成分类器的设计。而分类决策在识别过程中起作用,对待识别的样品进行分类决策。

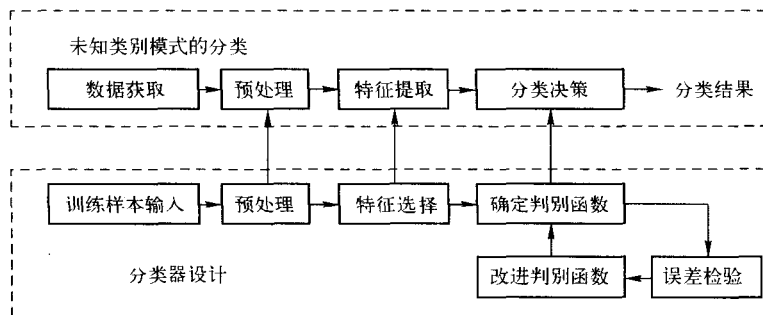


图 1-1 模式识别的过程

模式识别系统组成单元功能如下。

(1) 数据获取

用计算机可以运算的符号来表示所研究的对象,一般获取的数据类型有以下几种。

- ① 二维图像:文字、指纹、地图、照片等。
- ② 一维波形:脑电图、心电图、季节震动波形等。
- ③ 物理参量和逻辑值:体温、化验数据、参量正常与否的描述。

(2) 预处理

对输入测量仪器或其他因素所造成的退化现象进行复原、去噪声,提取有用信息。

(3) 特征提取和选择

对原始数据进行变换,得到最能反映分类本质的特征。将维数较高的测量空间(原始数据组成的空间)转变为维数较低的特征空间(分类识别赖以进行的空间)。

(4) 分类决策

在特征空间中用模式识别方法把被识别对象归为某一类别。

(5) 分类器设计

基本做法是在样品训练集基础上确定判别函数,改进判别函数和误差检验。

研究模式识别的主要目的是利用计算机进行模式识别,并对样本进行分类。执行模式识别的计算机系统称为模式识别系统。设计人员按需要设计模式识别系统,而该系统被用来执行模式分类的具体任务。

3. 统计模式识别研究的主要问题

统计模式识别主要研究的问题有:特征的选择与优化、分类判别、聚类判别。

(1) 特征的选择与优化

如何确定合适的特征空间是设计模式识别系统一个十分重要的问题,对特征空间进行优化有两种基本方法。一种是特征选择,如果所选用的特征空间能使同类物体分布具有紧致性,可以为分类器设计成功提供良好的基础;反之,如果不同类别的样品在该特征空间中混杂在一起,再好的设计方法也无法提高分类器的准确性。另一种是特征的组合优化,通过一种映射变换改造原特征空间,构造一个新的精简的特征空间。