

研究生教育“十二五”规划教材

# 模式识别原理及应用

余正涛 郭剑毅 毛存礼 等 编著

科学出版社

北京

## 内 容 简 介

本书系统介绍模式识别的基本理论和基本方法，主要内容包括统计模式识别、结构模式识别、模糊模式识别、神经网络模式识别等内容，加入当前应用广泛的隐马尔可夫模型、条件随机场模型、最大熵模型学习算法等内容，并通过大量实例讲述如何将所学理论知识运用到实际应用之中，对文本分类、文本聚类、语音识别、图像识别等应用做详细介绍，注重对主要知识内容的深入讨论，又突出新颖性。

本书可以作为计算机、自动化、通信工程及电子工程等专业高年级本科生和研究生的模式识别课程教材，也可供从事相关专业的教学、科研和工程技术人员参考。

### 图书在版编目 (CIP) 数据

模式识别原理及应用 /余正涛等编著. —北京：科学出版社，2014.6

ISBN 978-7-03-040570-8

I. ①模… II. ①余… III. ①模式识别 IV. ①O235

中国版本图书馆 CIP 数据核字 (2014) 第 094551 号

责任编辑：杨 岭 孟 锐 / 责任校对：郑金红

责任印制：余少力 / 封面设计：墨创文化

科 学 出 版 社 出 版

北京东黄城根北街 16 号

邮 政 编 码：100717

<http://www.sciencep.com>

成都创新包装印刷厂印刷

科学出版社发行 各地新华书店经销

\*

2014 年 6 月第 一 版 开本：B5 (720 × 1000)

2014 年 6 月第一次印刷 印张：26 1/2

字数：590 000

定 价：55.00 元

(如有印装质量问题，我社负责调换)

## 作者简介

余正涛，男，1970年出生，博士、教授、博士研究生导师，昆明理工大学信息工程与自动化学院院长，昆明理工大学智能信息处理重点实验室主任，昆明理工大学智能信息处理创新团队首席教授。2005年博士毕业于北京理工大学计算机应用技术专业，2008年11月至2009年12月在美国普渡大学做访问学者，主要从事自然语言处理、信息检索及机器学习方面的研究。入选首批中央组织部国家高层次人才特支计划“万人计划（科技创新领军人才）”，入选首批科学技术部创新人才推进计划“中青年科技创新领军人才”，云南省中青年学术技术带头人、国家自然科学基金项目通信评议人、中国科技奖励评审专家、中国计算机学会高级会员、中国中文信息学会理事、中国自动化学会理事、中国中文信息学会信息检索专委会委员、中国计算机学会协同计算专委会委员、中国计算机学会中文信息技术专委会委员、中国计算机学会互联网专委会委员、中国人工智能学会机器学习专委会委员。2008年享受云南省政府特殊津贴，主持国家自然科学基金、国家中小企业创新基金、教育部自然科学基金、云南省自然科学基金重点项目、云南省“九五”攻关项目、云南省教育厅基金及横向合作项目近40余项。发表学术论文150余篇，被SCI、EI收录80余篇，第一授权人授权国家发明专利2项，受理国家发明专利5项，授权软件著作权65项。以第一获奖人获得云南省科技进步奖一等奖、云南省自然科学奖二等奖、中国技术市场协会金桥奖、云南省科技进步奖三等奖各1项。

郭剑毅，女，1964年出生，昆明理工大学信息工程与自动化学院教授、硕士研究生导师。1990年硕士毕业于西安交通大学信息与控制工程系。中国计算机学会会员、中国中文信息学会会员。从事模式识别、决策分析与决策支持等教学工作多年。主持并参与了国家自然科学基金、云南省信息专项基金、云南省自然科学基金、云南省教育厅基金及昆明理工大学校青年基金等多项项目。第一授权人授权国家发明专利1项，授权国家软件著作权30余项；以第一作者发表论文60余篇，被SCI、EI收录20余篇；并获得云南省科技进步奖一等奖1项、云南省自然科学奖二等奖1项。主要从事自然语言处理、信息抽取和机器学习方面的研究。

# 序

模式识别是信息学科和人工智能的一个极其重要的分支，是一门理论与应用并重的技术科学。它以概率统计分析、模糊数学、神经网络、机器学习理论、句法结构和计算机信息处理等技术为基础，研究客观模式的机器分类(聚类)算法及其实现。模式识别的原理和方法已在很多地方得到了成功地应用，在医学、工程、军事、公安等众多领域应用十分广泛，是信息与计算及其相关专业进行科学研究的基础，但还有很多新问题需要应用模式识别技术去解决。

这门课的教学目的是让学生了解模式识别的基本概念、基本原理，掌握模式识别的基本分析方法和算法。通过对模式识别的基本理论和方法、运用实例的学习，使学生了解和掌握模式识别的基本理论与方法，具有初步设计、实现模式识别中比较简单的分类器算法的能力，培养学生利用模式识别方法、运用技能解决本专业及相关领域实际问题的能力，为将来继续深入学习或进行科学研究打下坚实的基础。

自然语言处理是当今的研究热点和难点，模式识别的一个重要应用之一就是自然语言处理。但目前的教材中介绍模式识别技术在自然语言处理方面的应用还不多。余正涛博士和他的合作者结合多年的教学实践和在自然语言处理等领域研究应用方面的研究成果，编写完成此书；书中也不乏相关技术在其他领域的应用内容，可以为相关研究提供很好的借鉴。

该书全面介绍统计模式识别、结构模式识别、模糊模式识别、神经网络模式识别、句法模式识别等基础理论，其应用部分穿插在每章之中，在该书的后几个章节中，详细介绍文本分类、文本聚类、语音识别和图像识别等应用。将理论与实际相结合，有利于读者加深对理论方法的理解，可使读者较系统地掌握模式识别的理论精髓和相关技术。书中给出的应用实例，为科研人员应用模式识别方法解决相关领域的实际问题提供了具体思路和方法。

纵观全书，作者系统地梳理模式识别技术的理论与方法及其最新成果和研究进展，继承传统模式识别理论和方法基础，探讨模式识别的相关理论、热点领域和主要技术方法，增添其在文本分类、信息检索、信息抽取等领域的应用。最后，通过文本分类、聚类、语音识别等系统的搭建和实验，对模式识别的典型方法给

予示范。全书结构严谨、内容丰富，具有一定前沿性和实用性，对模式识别在不同领域的应用有很好的参考价值。

章 肖

2013年10月12日

# 前　　言

自 20 世纪 60 年代以来，模式识别得到了迅速发展，并取得了丰富的理论成果，其应用领域也已扩展到了文本分类、语音识别、图像识别、视频识别、信息检索与数据挖掘等领域。由于模式识别理论具有重要的学术价值和广泛的应用领域，因而越来越多的人认识到模式识别课程的重要性，相关领域的科研工作者也投入了很高的学习热情。为了给在校本科生和研究生提供一本内容较新、论述较系统的有关模式识别的教材，也为了给相关领域的科研人员提供一本内容涵盖面广、具有一定前沿性和实用性的参考书，我们编写了这本书。

本书以基础理论教学为主，同时穿插实际应用来加深对基础理论的理解。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。书中不断引入和介绍学科最新的成果，增加模式识别在语音识别、信息检索等领域的应用。本书具有以下特点。

(1) 内容广泛：本书系统阐述模式识别领域的基础知识及经典方法，对经实践证明具有重要现实意义的新理论、新方法、新技术也进行介绍。包括统计模式识别、句法模式识别、模糊模式识别、神经网络技术和统计语言模型与信息检索模型，以及目前应用广泛的机器学习模型：隐马尔可夫模型(HMM)、最大熵模型(ME)和条件随机场模型(CRFs)等。

(2) 结构清晰合理：本书内容以学习目标、内容讲解、小结、习题与问题思考、参考文献为主线，按由浅入深、先易后难、先理论后应用、先传统后前沿来安排，有益于读者对各种理论、方法的理解。

(3) 选材考究精细：模式识别理论、方法、技术纷繁众多，而且新的理论和方法还在不断地产生。本书在众多的知识中选取基础理论、经典学习方法、典型应用等重要内容。

(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) 注重实践与应用：本书专门安排 4 章内容介绍模式识别的典型应用及实现过程，结合实例讲述模式识别的理论与方法，从而对读者理解模式识别理论与方法有很好的指导作用。

本书是余正涛、郭剑毅、毛存礼、线岩团、李华锋、王红斌、汤宏颖、张亚飞、高盛祥等教师共同努力的结果，其中第2章由余正涛编写，并负责制定全书大纲和编写组织，第1章、第10章、第14章由郭剑毅编写，并负责全书统稿，第13章、第20章由毛存礼编写，第16章、第18章、第19章由线岩团编写，第3章、第5章、第7章由李华锋编写，第8章、第9章、第12章由王红斌编写，第6章、第11章由汤宏颖编写，第4章、第15章由张亚飞编写，第17章由高盛祥编写。同时感谢参与本教材编写工作的研究生，他们是洪旭东、吴则建、陈方琼、田维、宋海霞、潘霄、石林宾、陈鹏、康潮明、张优敏、赵君、李真、魏斯超、潘清清、于海涛、潘华山、王炎冰等。

衷心感谢科学出版社的编辑，是他们的辛勤劳动，使本书得以顺利出版。最后，本书参考了国内外许多同行的论文、著作，引用了其中的观点、数据与结论，在此一并表示谢意。

由于作者学识有限，书中不足之处在所难免，敬请批评、指正。

编 者

2013年10月

# 目 录

<b>第 1 章 模式识别概论 .....</b>	1
1.1 概述 .....	1
1.2 模式识别的发展历史 .....	3
1.3 模式识别与其他学科的关系 .....	5
1.4 模式识别的基本方法 .....	6
1.5 模式识别的应用 .....	8
本章小结 .....	12
习题与思考题 .....	15
参考文献 .....	15
<b>第 2 章 模式识别的基本概念 .....</b>	17
2.1 概述 .....	17
2.2 基本概念 .....	17
2.3 模式识别系统 .....	20
2.4 模式识别的一些基本问题 .....	21
2.5 相关数学概念 .....	27
本章小结 .....	30
习题与思考题 .....	31
参考文献 .....	31
<b>第 3 章 模式识别的判别函数 .....</b>	32
3.1 概述 .....	32
3.2 线性判别函数的基本概念 .....	32
3.3 线性判别函数的判定面 .....	33
3.4 非线性判别函数 .....	40
3.5 广义线性判别函数 .....	43
3.6 线性分类器的设计 .....	50
本章小结 .....	53
习题与思考题 .....	53
参考文献 .....	54

---

7.4 聚类准则函数 .....	129
7.5 聚类算法 .....	134
7.6 聚类分析在虚假评论检测中的应用 .....	146
本章小结 .....	150
习题与思考题 .....	151
参考文献 .....	152
<b>第 8 章 句法模式识别 .....</b>	<b>153</b>
8.1 概述 .....	153
8.2 形式语言概述 .....	155
8.3 基元提取和文法推断 .....	159
8.4 句法分析 .....	161
8.5 自动机理论 .....	164
8.6 自动机理论在语音识别中的应用 .....	168
本章小结 .....	172
习题与思考题 .....	173
参考文献 .....	174
<b>第 9 章 模糊模式识别 .....</b>	<b>175</b>
9.1 概述 .....	175
9.2 模糊模式识别的基本概念 .....	175
9.3 直接模糊模式识别法 .....	177
9.4 间接模糊模式识别法 .....	181
9.5 模糊聚类 .....	183
9.6 模糊模式识别在大气质量评定中的应用 .....	190
本章小结 .....	193
习题与思考题 .....	193
参考文献 .....	195
<b>第 10 章 决策树 .....</b>	<b>196</b>
10.1 概述 .....	196
10.2 决策树学习 .....	197
10.3 CLS 学习算法 .....	199
10.4 ID3 学习算法 .....	202
10.5 决策树的剪枝技术 .....	210
10.6 决策树的评价 .....	215

10.7 决策树算法的优化 .....	215
10.8 决策树的应用 .....	216
10.8.1 决策树在文本分类中的应用 .....	216
10.8.2 基于决策树的个人住房贷款信用风险评估模型 .....	220
本章小结 .....	226
习题与思考题 .....	227
参考文献 .....	229
<b>第 11 章 人工神经网络 .....</b>	<b>230</b>
11.1 概述 .....	230
11.2 神经元 .....	230
11.3 人工神经网络拓扑结构 .....	233
11.4 人工神经网络学习方法及规则 .....	233
11.5 前馈神经网络及其主要算法 .....	235
11.6 Hopfield 网络 .....	241
11.7 自组织神经网络 .....	244
11.8 人工神经网络的应用 .....	246
本章小结 .....	248
习题与思考题 .....	249
参考文献 .....	250
<b>第 12 章 隐马尔可夫模型 .....</b>	<b>251</b>
12.1 概述 .....	251
12.2 隐马尔可夫模型的概念 .....	251
12.3 隐马尔可夫模型的三个基本问题及解决办法 .....	256
12.4 隐马尔可夫模型在中文旅游景点识别中的应用 .....	265
本章小结 .....	269
习题与思考题 .....	270
参考文献 .....	271
<b>第 13 章 最大熵模型 .....</b>	<b>272</b>
13.1 概述 .....	272
13.2 熵及最大熵 .....	273
13.3 最大熵模型 .....	280
13.4 最大熵在自然语言处理中的应用 .....	286
本章小结 .....	290

---

习题与思考题 .....	290
参考文献 .....	291
<b>第 14 章 条件随机场 .....</b>	<b>292</b>
14.1 概述 .....	292
14.2 概率图模型 .....	292
14.3 条件随机场简介 .....	298
14.4 势函数 .....	299
14.5 参数估计与训练 .....	300
14.6 参数估计的优化 .....	305
14.7 条件随机场在旅游领域命名实体识别中的应用 .....	307
本章小结 .....	312
习题与思考题 .....	313
参考文献 .....	313
<b>第 15 章 统计学习理论及支持向量机 .....</b>	<b>315</b>
15.1 概述 .....	315
15.2 机器学习的基本问题和方法 .....	316
15.3 统计学习理论 .....	318
15.4 支持向量机 .....	322
15.5 支持向量机的分类与回归 .....	329
15.6 基于支持向量机的汉语问句分类 .....	338
本章小结 .....	341
习题与思考题 .....	342
参考文献 .....	342
<b>第 16 章 统计语言模型及信息检索 .....</b>	<b>344</b>
16.1 概述 .....	344
16.2 统计语言模型 .....	344
16.3 信息检索 .....	353
16.4 统计语言模型在拼音输入法中的应用 .....	367
本章小结 .....	370
习题与思考题 .....	371
参考文献 .....	371
<b>第 17 章 基于 SVM 的中文文本分类 .....</b>	<b>373</b>
17.1 概述 .....	373

可以不断地根据周围的景物，判断他是否能达到目的地，这实际也是不断地在作“正确”和“不正确”的分类判断。人脑的这种思维能力就构成了“模式”的概念。

在狭义上，模式识别这门课，仅局限于研究具体的客观事物，以及如何用计算机或机器来进行自动识别。因此，狭义上说，模式是对客观事物的一种定量的或结构的描述，而具有某些共同特征的模式的集合被称为模式类，其中个别具体的模式往往称为样本。模式识别就是样本到类别的映射。因为模式识别是使机器能自动地(或有人工少量干预)把待识别模式(被识对象)分配到各自的模式类中去，因此，常常又把模式识别称作模式分类或机器识别。

## 1.2 模式识别的发展历史

模式识别诞生于 20 世纪 20 年代，当时已有用光学和机械手段实现模式识别的例子，如在 1929 年 Gustav Tauschek 就在德国获得了光学字符识别的专利，发明了阅读机，能够阅读 0~9 的数字<sup>[1]</sup>。随后在 30 年代 Fisher 提出统计分类理论<sup>[1]</sup>，奠定了统计模式识别的基础。20 世纪 40 年代计算机的出现，增加了对模式识别实际应用的需求，也推动了模式识别理论的发展。1957 年 IBM 的 C.K.Chow 将统计决策方法用于字符识别<sup>[2]</sup>；同样在 50 年代美国语言学家、转换生成语法的创始人 Noam Chomsky 提出了形式语言理论，美籍华人傅京荪(K. S. Fu)提出句法结构模式识别<sup>[3]</sup>；60 年代美国控制论学者 L.A.Zadeh 提出了模糊集(fuzzy set)概念，建立了模糊集理论，模糊模式识别理论得到了较广泛的应用<sup>[4]</sup>。然而，“模式识别”这个词被广泛使用并形成一个领域则是在 20 世纪 60 年代以后。1966 年由 IBM 组织在波多黎各召开了第一次以“模式识别”为题的学术会议。Nagy<sup>[5,6]</sup> 和 Kanal<sup>[7]</sup> 分别介绍了 1968 年以前和 1968~1974 年的研究进展。70 年代几本很有影响的模式识别教材(如 Fukunaga<sup>[8]</sup>, Duda 和 Hart<sup>[9]</sup>)的相继出版和 1972 年 IEEE 发起的第一届国际模式识别大会(ICPR)的召开，标志着模式识别领域的形成；同时，国际模式识别协会(IAPR)在 1974 年的第二届国际模式识别大会上开始筹建，模式识别的国际会议“ICPR”在 1978 年的第四届大会上正式成立。1977 年 IEEE 的计算机学会成立了模式分析与机器智能(PAMI)委员会，每两年召开一次模式识别与图像处理学术会议。国内方面，由模式识别国家重点实验室主办，中国自动化学会、中国图象图形学学会协办的第一届全国模式识别学术会议(CCP)于 2007 年 12 月 11 至 12 日在北京召开；第一届亚洲模式识别会议(ACPR 2011)也于 2011 年 11 月 28 日至 30 日在中国首都北京举行。这些会议的召开，旨在进一步促进模式识别研究的快速发展，加强国内外同行间的学术交流与合作，

## 1.3 模式识别与其他学科的关系

模式识别是一个综合性的交叉学科，它与很多学科都有联系，特别是与数学和人工智能的联系十分密切。模式识别的理论基础是统计学和一些近代数学方法，除此之外，模式识别还涉及计算机科学、心理学、语言学、工程学、控制论以及生物医学等众多领域。

### 1.3.1 模式识别与数学

模式识别是一门与数学结合非常紧密的科学，所应用到的数学知识非常多，最基本的便是概率论和数理统计。模式识别技术到处都充满了概率和统计的思想，我们经常所说的识别率，其实就是概率的表达：在大数据量(严格地说应当是数据量无穷大)测试中识别成功的概率；还有常用的贝叶斯决策分类器便是运用了概率公式。模式识别还用到了线性代数，因为运用线性代数可以较方便表达具有多特征的事物，我们一般会用向量来表达一个事物的特征，对于向量的计算是一定会用到线性代数的知识的。还有一个较为高层次的数学知识是泛函分析，泛函分析是研究无限维线性空间上的泛函数和算子理论，SVM 便是以泛函分析中的理论为基础的，SVM 技术还运用到了最优化理论。

### 1.3.2 模式识别与人工智能

早期的模式识别研究是与人工智能和机器学习密不可分的，如 Rosenblatt 的感知机<sup>[11]</sup>和 Nilsson 的学习机<sup>[12]</sup>就与这三个领域密切相关。后来，由于人工智能更关心符号信息和知识的推理，而模式识别更关心感知信息的处理，二者逐渐分离形成了不同的研究领域。介于模式识别和人工智能之间的机器学习在 20 世纪 80 年代以前也偏重于符号学习，后来人工神经网络重新受到重视，统计学习逐渐成为主流，与模式识别中的学习问题渐趋重合，重新拉近了模式识别与人工智能的距离。模式识别与机器学习的方法也被广泛用于感知信号以外的数据分析问题(如文本分析、商业数据分析、基因表达数据分析等)，形成了数据挖掘领域。

模式分类是模式识别的主要任务和核心研究内容。分类器设计是在训练样本集合上进行优化(如使每一类样本的表达误差最小或使不同类别样本的分类误差最小)的过程，也就是一个机器学习过程。由于模式识别的对象是存在于感知信号中的物体和现象，它研究的内容还包括对图像和视频信号的处理、分割、形状和运动分析等，以及面向应用(如文字识别、语音识别、生物认证、医学图像分析、遥感图像分析等)的方法和系统研究。

现在研究的大部分都是统计模式识别的方法，而且在这其中研究比较集中的

用字符串或图来表示；然后运用形式语言理论进行句法分析，依据其是否符合某一类的文法而决定其类别。基元组合成模式的规则，由所谓语法来指定。一旦基元被鉴别，识别过程可通过句法分析进行，即分析给定的模式语句是否符合指定的语法，满足某类语法的即被分入该类。这种方法的优点是适合结构性强的模式，缺点是抗噪声能力差、计算复杂度高。

#### 1.4.4 模糊模式识别

模糊模式识别(fuzzy pattern recognition)是基于模糊数学的识别方法。现实世界中存在许多界限不分明、难以精确描述的事物或现象，而模糊数学则可以用数学的方法研究和处理这类具有“模糊性”的事物或现象。模糊数学的出现使得人们可以模拟人类神经系统的活动，描述模式属于某类的程度，因此，模糊数学在模式分类中得到了很好的应用。目前，模糊模式识别方法较多，比如模糊近邻、模糊最小最大神经网络等。该类方法的有效性主要在于对象类的隶属函数是否良好。

#### 1.4.5 神经网络模式识别

神经网络模式识别(neural network pattern recognition)利用神经元网络中出现的神经计算模式进行。大部分神经元网络都有某种训练规则，如基于现有模式调节连接权重。换句话说，神经元网络直接对例子进行学习，得出其结构特征进行推广，就像孩子从狗的例子中认识狗一样。人工神经元网络可以超越传统基于计算机的模式分类系统的能力。人们可以利用计算机或神经元网络进行模式分类。计算机利用传统的数学算法来检测给定的模式是否与现有模式相匹配，这是一个简单易懂的方法。但是，该方法只能进行是或非的判断，且不允许模式有噪声。神经元网络允许模式有噪声，而且如果训练得当，神经元网络会对未知模式的类别做出正确的响应。虽然神经元网络不能创造奇迹，但是如果采用合适的结构，对好的数据进行正确的训练，神经元网络都可以给出令人惊异的结果。比如，BP 神经网络直接从观测数据(训练样本)学习，非常简便有效，因而得到了广泛应用。

在上述分类方法中，统计模式识别与结构模式识别是模式识别中的经典性和基础性技术；而模糊识别与神经网络识别则是最近发展起来的新方法，是信息科学与人工智能的重要组成部分。另外，上述几种分类方法各有自己的特点与应用范围，它们不能相互取代，只能相互共存、促进、借鉴、渗透与融合。方法的选择取决于问题的性质。一个好的分类方法可能综合利用了上述各类识别方法的观点、概念和技术而形成。总之，模式识别作为一门学科，尽管它已经经历了几十年的历史，建立了丰富的理论体系，但由于问题的复杂性以及其应用领域范围的

为人工键盘输入和机器自动识别输入两种。其中人工键入速度慢而且劳动强度大；自动输入又分为汉字识别输入及语音识别输入。从识别技术的难度来说，手写体识别的难度高于印刷体识别，而在手写体识别中，脱机手写体的难度又远远超过了联机手写体识别。到目前为止，除了脱机手写体数字的识别已有实际应用外，汉字等文字的脱机手写体识别还处在实验阶段。

### 1.5.3 语音识别

模式识别中的一个重要应用是语音识别，其目的就是让计算机能听懂人说的话。语音识别技术的应用包括语音拨号、语音导航、室内设备控制、语音文档检索、简单的听写数据录入、音乐搜索等。语音识别技术是涉及声学、语言学、数字信号处理、统计模式识别、概率论和信息论、发声机理和听觉机理、人工智能等多学科技术的一项综合性技术。近年来，在生物识别技术领域中，声纹识别技术以其独特的方便性、经济性和准确性等优势受到世人瞩目，并日益成为人们日常生活和工作中重要且普及的安全验证方式。目前，主流的大词汇量语音识别系统多采用统计模式识别技术。随着智能移动设备的普及，语音交互作为一种新型的人机交互方式，正越来越引起整个IT业界的重视。语音交互是一个认知过程，不能与语言的语法、语义和语用结构割裂开来，所以语音识别的主要问题不是识别，更多的是语义分析和理解、上下文分析以及用户体验等。语音识别技术与自然语言处理技术如机器翻译及语音合成技术相结合，可以构建出更加复杂的应用，例如语音到语音的翻译等，正逐步成为信息技术中人机接口(human computer interface, HCI)的关键技术。

### 1.5.4 计算机辅助医学诊断

计算机辅助医学诊断是模式识别的重要应用之一，是指通过影像学、医学图像处理技术以及其他可能的生理、生化手段，结合计算机的分析计算，提高诊断的准确率。计算机辅助诊断的系统需求来源于如下事实：医疗数据较难解释并且解释结果多依赖于医生的经验。通常医学影像学中计算机辅助诊断分为三步：第一步是把病变从正常结构中提取出来；第二步是图像特征的量化；第三步是对数据进行处理并得出结论。因为计算机可以全面利用影像信息进行精确的定量计算，去除人的主观性，避免因个人知识和经验的差异而引起的“千差万别”的诊断结果，所以它的结果是不含糊的、确定的，它使诊断变得更为准确、科学。计算机辅助诊断已经应用于实际，主要研究各种医疗数据，如X射线、计算机断层图、超声波图、心电图和脑电图。目前在癌细胞检测、X射线照片分析、血液化验、染色体分析、心电图诊断和脑电图诊断等方面，模式识别已取得了成效。计算机辅助诊断技术虽然在国内已有一定的开展，但其发展缓慢，尚有许多

### 1.5.8 文本挖掘

文本挖掘(text mining)是指为了发现知识，从大规模文本库中抽取隐含的、以前未知的、潜在有用的知识(包括概念、模式、规则、规律、约束等形式)。文本挖掘的过程一般包括文本数据预处理、文本信息提取和索引、文本知识挖掘及知识后处理等步骤。数据预处理包括数据清洗(如去噪、去重)、数据选择(选择合适的、面向特定领域的文本数据)和文本切分(如中文分词、段落切分)等。数据预处理后，必须提取中文文本的特征信息，包括关键词提取、术语提取、基于模板的信息抽取和基于专业词典的概念转换等操作。经过中文文本特征提取操作后，中文文本数据转换为中文文本信息，在文本信息的基础上进行知识挖掘，包括文本自动摘要、文本聚类、关联规则抽取和语义关系挖掘等。由于知识挖掘得到的结果可能不一致、不新颖、不符合构建本体基本要素的形式要求，因此需要对文本知识进行必要的后处理，包括知识的评价与取舍、知识的规范形式化表达等。文本挖掘是搜索引擎的重要技术，已经成为信息检索、数据挖掘、机器学习、模式识别、统计以及计算语言学等学科中的重要领域。

### 1.5.9 信息检索

信息检索(information retrieval)是指搜索信息的科学，根据用户的查询要求，从信息数据库、网页、Word文档、PDF文档、图像，甚至视频和音频文件等信息源中检索出与之相关的信息资料。在信息检索领域，英语信息检索的发展较为迅速。英语信息检索系统可以利用向量空间表示检索信息内容，并将自然语言处理应用于信息检索，大大提高了信息查询的准确性。中文语词之间由于没有空格，因此在索引前需要进行语词切分；另一方面，与英语相比，汉语句法分析和语义理解更为困难，造成中文信息检索的发展较为缓慢。信息检索过程就是一个模式识别的过程，首先对信息需求进行主题分析，形成能代表检索需求的概念，并将这些概念转换成信息检索语言的词语，然后在检索工具中进行匹配运算，从而找到所需的信息。因此，信息检索的实质是一个匹配过程，也就是信息用户需求的主题概念或检索表达式同一定信息系统的系统语言相匹配的过程，如果两者匹配，则所需信息被检中，否则检索失败。早期的信息检索模型虽然构造方法简单，但精度较低，难以获得用户满意的检索结果。对此，近年来国外有学者提出把排序学习(learning to rank)应用到检索模型的构造上，以期获得更精确的检索结果。信息检索的核心就是搜索引擎，互联网搜索与统计机器学习密不可分。搜索引擎能够通过机器学习有效地将海量数据联系、组织、利用起来，在大规模的分布式计算平台上为用户提供服务。

网络法和模糊识别法等。上述几种分类方法各有自己的特点与应用范围，它们不能相互取代，只能相互共存、促进、借鉴、渗透与融合。方法的选择取决于问题的性质。

统计模式识别的理论体系已经相当完善，方法也很多，如贝叶斯分类器、神经网络、SVM 法、Fisher 方法、KNN 法等。在统计方法中，特征抽取占有重要的地位，但尚无通用的理论指导，只能通过分析具体识别对象决定选取何种特征。

结构模式识别称为句法模式识别。在多数情况下，可以有效地用形式语言理论中的文法表示模式的结构信息，因此也常称为句法模式识别。结构模式识别主要立足于分析模式的结构信息。该技术将对象分解为若干个基本单元，即基元；用这些基元以及它们的结构关系来描述对象，基元以及这些基元的结构关系可以用字符串或图来表示；然后运用形式语言理论进行句法分析，依据其是否符合某一类的文法而决定其类别。基元组合成模式的规则，由所谓语法来指定。一旦基元被鉴别，识别过程可通过句法分析进行，即分析给定的模式语句是否符合指定的语法，满足某类语法的即被分入该类。这种方法的优点是适合结构性强的模式，缺点是抗噪声能力差、计算复杂度高。

模糊模式识别是基于模糊数学的识别方法。现实世界中存在许多界限不分明、难以精确描述的事物或现象，而模糊数学则可以用数学的方法研究和处理这类具有“模糊性”的事物或现象。模糊数学的出现使得人们可以模拟人类神经系统的活动，描述模式属于某类的程度，因此，模糊数学在模式分类中得到了很好的应用。目前，模糊模式识别方法较多，比如模糊-近邻、模糊最小最大神经网络等。该类方法的有效性主要在于对象类的隶属函数是否良好。

神经网络方法与统计方法相比具有不依赖概率模型、参数自学习、泛化性能良好等优点，至今仍在模式识别中广泛应用。然而，神经网络的设计和实现依赖于经验，泛化性能不能确保最优。

20 世纪 90 年代 SVM 的提出吸引了模式识别界对统计学习理论和核方法的极大兴趣。与神经网络相比，SVM 的优点是通过优化一个泛化误差界限自动确定一个最优的分类器结构，从而具有更好的泛化性能。而核函数的引入使很多传统的统计方法从线性空间推广到高维非线性空间，提高了表示和判别能力。

## 5. 模式识别已经在哪些领域成功应用？举例分析说明。

模式识别研究在近几年来取得了令人瞩目的成就，一批批研究成果在越来越多的领域得到了广泛应用和推广。如我们熟知的文献分类、财政预测、信件分拣、文字和语音识别、遥感图片的机器判读、生物特征识别(指纹、虹膜、人脸等识别)、生物医学的细胞或组织分析、系统的故障诊断、具有视觉的机器人、汽车自动驾驶系统、信息分类与检索、网络入侵检测、工业产品检测等领域，并且正在扩展