

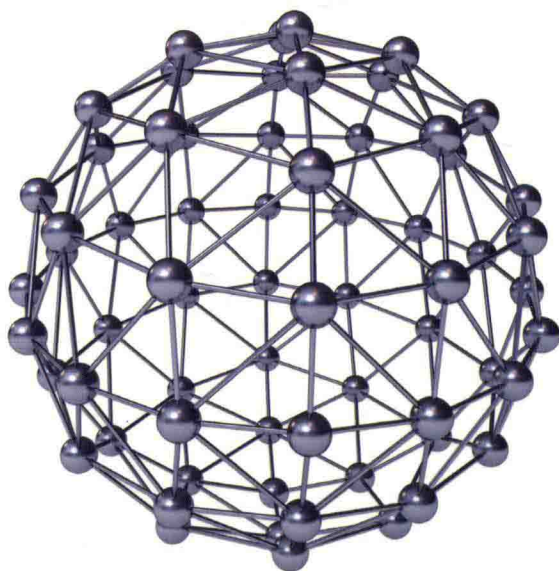


教育部高等学校计算机类专业教学指导委员会-华为ICT产学合作项目  
数据科学与大数据技术专业系列规划教材

华为信息与网络  
技术学院指定教材

# 机器学习

赵卫东 董亮◎编著



系统、完整的数据科学与大数据技术专业解决方案

名校名师打造大数据领域精品力作

强调基本概念和机器学习算法

兼顾机器学习经典内容，突出深度学习前沿

 中国工信出版集团

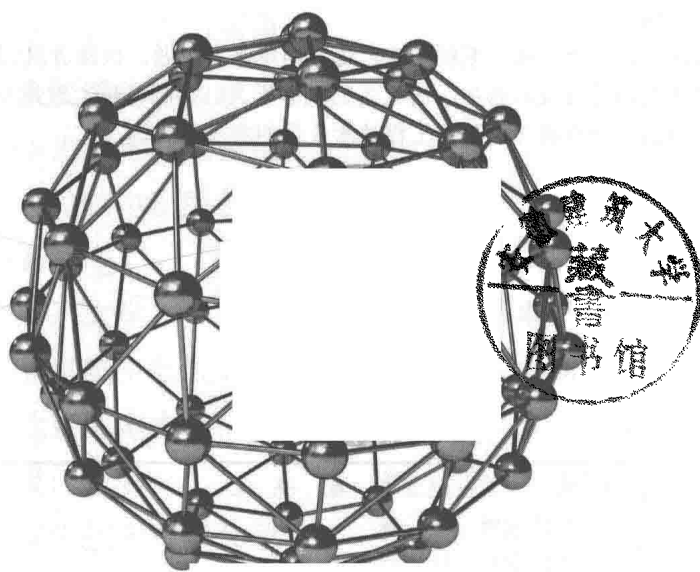
 人民邮电出版社  
POSTS & TELECOM PRESS

校计算机类专业教学指导委员会-华为ICT产学合作项目  
与大数据技术专业系列规划教材

华为信息与网络  
技术学院指定教材

# 机器学习

赵卫东 董亮◎编著



人民邮电出版社  
北京

## 图书在版编目 (C I P) 数据

机器学习 / 赵卫东, 董亮编著. — 北京: 人民邮电出版社, 2018. 8  
数据科学与大数据技术专业系列规划教材  
ISBN 978-7-115-48300-3

I. ①机… II. ①赵… ②董… III. ①机器学习—教材 IV. ①TP181

中国版本图书馆CIP数据核字(2018)第163066号

## 内 容 提 要

机器学习是人工智能的重要技术基础,涉及的内容十分广泛。本书涵盖了机器学习的基础知识,主要包括机器学习的概述、统计学基础、分类、聚类、神经网络、贝叶斯网络、支持向量机、进化计算、文本分析等经典的机器学习基础知识,还包括用于大数据机器学习的分布式机器学习算法、深度学习等高级内容。此外,本书还介绍了机器学习的热门应用领域推荐系统,并给出了华为机器学习平台上的实验。

本书深入浅出、内容全面、案例丰富,每章结尾都有习题,供读者巩固所学知识。

本书适合作为高等院校本科生、研究生的机器学习、数据分析、数据挖掘等课程的教材,也可作为对机器学习感兴趣的研究人员和工程技术人员的参考资料。

---

◆ 编 著 赵卫东 董 亮

责任编辑 张 斌

责任印制 沈 蓉 彭志环

◆ 人民邮电出版社出版发行 北京市丰台区成寿寺路 11 号

邮编 100164 电子邮件 315@ptpress.com.cn

网址 <http://www.ptpress.com.cn>

三河市潮河印业有限公司印刷

◆ 开本: 787×1092 1/16

印张: 23.25

2018年8月第1版

字数: 611千字

2018年8月河北第1次印刷

---

定价: 59.80 元

读者服务热线: (010)81055256 印装质量热线: (010)81055316

反盗版热线: (010)81055315

广告经营许可证: 京东工商广登字 20170147 号

教育部高等学校计算机类专业教学指导委员会-华为 ICT 产学合作项目  
数据科学与大数据技术专业系列规划教材

## 编 委 会

- 主任 陈 钟 北京大学  
副主任 杜小勇 中国人民大学  
周傲英 华东师范大学  
马殿富 北京航空航天大学  
李战怀 西北工业大学  
冯宝帅 华为技术有限公司  
张立科 人民邮电出版社  
秘书长 王 翔 华为技术有限公司  
戴思俊 人民邮电出版社

委 员 (按姓名拼音排序)

- |     |          |     |         |
|-----|----------|-----|---------|
| 崔立真 | 山东大学     | 段立新 | 电子科技大学  |
| 高小鹏 | 北京航空航天大学 | 桂劲松 | 中南大学    |
| 侯 宾 | 北京邮电大学   | 黄 岚 | 吉林大学    |
| 林子雨 | 厦门大学     | 刘 博 | 人民邮电出版社 |
| 刘耀林 | 华为技术有限公司 | 乔亚男 | 西安交通大学  |
| 沈 刚 | 华中科技大学   | 石胜飞 | 哈尔滨工业大学 |
| 嵩 天 | 北京理工大学   | 唐 卓 | 湖南大学    |
| 汪 卫 | 复旦大学     | 王 伟 | 同济大学    |
| 王宏志 | 哈尔滨工业大学  | 王建民 | 清华大学    |
| 王兴伟 | 东北大学     | 薛志东 | 华中科技大学  |
| 印 鉴 | 中山大学     | 袁晓如 | 北京大学    |
| 张志峰 | 华为技术有限公司 | 赵卫东 | 复旦大学    |
| 邹北骥 | 中南大学     | 邹文波 | 人民邮电出版社 |

毫无疑问，我们正处在一个新时代。新一轮科技革命和产业变革正在加速推进，技术创新日益成为重塑经济发展模式和促进经济增长的重要驱动力，而“大数据”无疑是第一核心推动力。

当前，发展大数据已经成为国家战略，大数据在引领经济社会发展中的新引擎作用更加突显。大数据重塑了传统产业的结构和形态，催生了众多的新产业、新业态、新模式，推动了共享经济的蓬勃发展，也给我们的衣食住行带来根本改变。同时，大数据是带动国家竞争力整体跃升和跨越式发展的巨大推动力，已成为全球科技和产业竞争的重要制高点。可以大胆预测，未来，大数据将会进一步激起全球科技和产业发展浪潮，进一步渗透到我们国计民生的各个领域，其发展扩张势不可挡。可以说，我们处在一个“大数据”时代。

大数据不仅仅是单一的技术发展领域和战略新兴产业，它还涉及科技、社会、伦理等诸多方面。发展大数据是一个复杂的系统工程，需要科技界、教育界和产业界等社会各界的广泛参与和通力合作，需要我们以更加开放的心态，以进步发展的理念，积极主动适应大数据时代所带来的深刻变革。总体而言，从全面协调可持续健康发展的角度，推动大数据发展需要注重以下五个方面的辩证统一和统筹兼顾。

一是要注重“长与短结合”。所谓“长”就是要目标长远，要注重制定大数据发展的顶层设计和中长期发展规划，明确发展方向和总体目标；所谓“短”就是要着眼当前，注重短期收益，从实处着手，快速起效，并形成效益反哺的良性循环。

二是要注重“快与慢结合”。所谓“快”就是要注重发挥新一代信息技术产业爆炸性增长的特点，发展大数据要时不我待，以实际应用需求为牵引加快推进，力争快速占领大数据技术和产业制高点；所谓“慢”就是防止急功近利，欲速而不达，要注重夯实大数据发展的基础，着重积累发展大数据基础理论与核心共性关键技术，培养行业领域发展中的大数据思维，潜心培育大数据专业人才。

三是要注重“高与低结合”。所谓“高”就是要打造大数据创新发展高地，要结合国家重大战略需求和国民经济主战场核心需求，部署高端大数据公共服务平台，组织开展国家级大数据重大示范工程，提升国民经济重点领域和标志性行业的大数据技术水平和应用能力；所谓“低”就是要坚持“润物细无声”，推进大数据在各行各业和民生领域的广泛应用，推进大数据发展的广度和深度。

四是要注重“内与外结合”。所谓“内”就是要向内深度挖掘和深入研究大数据作为一门学科领域的深刻技术内涵，构建和完善大数据发展的完整理论体系和技术支撑体系；所谓“外”就是要加强开放创新，由于大数据涉及众多学科领域和产业行业门类，也涉及国家、社会、个人等诸多问题，因此，需要推动国际国内科技界、产业界的深入合作和各级政府广泛参与，共同研究制定标准规范，推动大数据与人工智能、云计算、物联网、网络安全等信息技术领域的协同发展，促进数据科学与计算机科学、基础科学和各种应用科学的深度融合。

五是要注重“开与闭结合”。所谓“开”就是要坚持开放共享，要鼓励打破现有体制机制障碍，推动政府建立完善开放共享的大数据平台，加强科研机构、企业间技术交流合作，推动大数据资源高效利用，打破数据壁垒，普惠数据服务，缩小数据鸿沟，破除数据孤岛；所谓“闭”就是要形成价值链生态闭环，充分发挥大数据发展中技术驱动与需求牵引的双引擎作用，积极运用市场机制，形成技术创新链、产业发展链和资金服务链协同发展的态势，构建大数据产业良性发展的闭环生态圈。

总之，推动大数据的创新发展，已经成为了新时代的新诉求。刚刚闭幕的党的十九大更是明确提出要推动大数据、人工智能等信息技术产业与实体经济深度融合，培育新增长点，为建设网络强国、数字中国、智慧社会形成新动能。这一指导思想为我们未来发展大数据技术和产业指明了前进方向，提供了根本遵循。

习近平总书记多次强调“人才是创新的根基”“创新驱动实质上是人才驱动”。绘制大数据发展的宏伟蓝图迫切需要创新人才培养体制机制的支撑。因此，需要把高端人才队伍建设作为大数据技术和产业发展的重中之重，需要进一步完善大数据教育体系，加强人才储备和梯队建设，将以大数据为代表的新兴产业发展对人才的创新性、实践性需求渗透融入人才培养各个环节，加快形成我国大数据人才高地。

国家有关部门“与时俱进，因时施策”。近期，国务院办公厅正式印发《关于深化产教融合的若干意见》，推进人才和人力资源供给侧结构性改革，以适应创新驱动发展战略的新形势、新任务、新要求。教育部高等学校计算机类专业教学指导委员会、华为公司和人民邮电出版社组织编写的《教育部高等学校计算机类专业教学指导委员会-华为 ICT 产学合作项目——数据科学与大数据技术专业系列规划教材》的出版发行，就是落实国务院文件精神，深化教育供给



侧结构性改革的积极探索和实践。它是国内第一套成专业课程体系规划的数据科学与大数据技术专业系列教材，作者均来自国内一流高校，且具有丰富的数据教学、科研、实践经验。它的出版发行，对完善大数据人才培养体系，加强人才储备和梯队建设，推进贯通大数据理论、方法、技术、产品与应用等的复合型人才培养，完善大数据领域学科布局，推动大数据领域学科建设具有重要意义。同时，本次产教融合的成功经验，对其他学科领域的人才培养也具有重要的参考价值。

我们有理由相信，在国家战略指引下，在社会各界的广泛参与和推动下，我国的大数据技术和产业发展一定会有光明的未来。

是为序。



中国科学院院士 郑志明

2018年4月16日





在 500 年前的大航海时代，哥伦布发现了新大陆，麦哲伦实现了环球航行，全球各大洲从此连接了起来，人类文明的进程得以推进。今天，在云计算、大数据、物联网、人工智能等新技术推动下，人类开启了智能时代。

面对这个以“万物感知、万物互联、万物智能”为特征的智能时代，“数字化转型”已是企业寻求突破和创新的必由之路，数字化带来的海量数据成为企业乃至整个社会最重要的核心资产。大数据已上升为国家战略，成为推动经济社会发展的新引擎，如何获取、存储、分析、应用这些大数据将是这个时代最热门的话题。

国家大数据战略和企业数字化转型成功的关键是培养多层次的大数据人才，然而，根据计世资讯的研究，2018 年中国大数据领域的人才缺口将超过 150 万人，人才短缺已成为制约产业发展的突出问题。

2018 年初，华为公司提出新的愿景与使命，即“把数字世界带入每个人、每个家庭、每个组织，构建万物互联的智能世界”，它承载了华为公司的历史使命和社会责任。华为企业 BG 将长期坚持“平台+生态”战略，协同生态伙伴，共同为行业客户打造云计算、大数据、物联网和传统 ICT 技术高度融合的数字化转型平台。

人才生态建设是支撑“平台+生态”战略的核心基石，是保持产业链活力和持续增长的根本，华为以 ICT 产业长期积累的技术、知识、经验和成功实践为基础，持续投入，构建 ICT 人才生态良性发展的使能平台，打造全球有影响力的 ICT 人才认证标准。面对未来人才的挑战，华为坚持与全球广大院校、伙伴加强合作，打造引领未来的 ICT 人才生态，助力行业数字化转型。

一套好的教材是人才培养的基础，也是教学质量的重要保障。本套教材的出版，是华为在大数据人才培养领域的重要举措，是华为集合产业与教育界的高端智力，全力奉献的结晶和成果。在此，让我对本套教材的各位作者表示由衷的感谢！此外，我们还要特别感谢教育部高等学校计算机类专业教学指导委员会副主任、北京大学陈钟教授以及秘书长、北京航空航天大学马殿富教授，没有你们的努力和推动，本套教材无法成型！

同学们、朋友们，翻过这篇序言，开启学习旅程，祝愿在大数据的海洋里，尽情展示你们的才华，实现你们的梦想！



华为公司董事、企业 BG 总裁 阎力大

2018 年 5 月



机器学习是人工智能的技术基础,伴随着人工智能几十年的发展,期间有过几次大起大落。作为机器学习的高级阶段,最近几年,深度学习算法在自然语言处理、语音识别、图像处理等领域的突破,使得机器学习成为计算机学科非常热门的一个方向。这也标志着机器学习已经彻底迈出实验室大门,走向实践,推动着人工智能向更高阶段发展。

与机器学习十分密切的概念有数据挖掘、大数据分析等,这些数据分析技术使用了一些机器学习的方法和算法,解决了企业应用的一些问题,辅助业务人员和管理人员做出更好的决策。几种技术相辅相成,共同促进了数据分析技术和人工智能的进步。

从早期的统计学习,发展到联结主义的神经网络,直到深度神经网络的过程中,机器学习的基础一直是足量的高质量数据。“互联网+”的热潮推动了大数据的产生以及处理大数据的软硬件技术的迅猛发展,为机器学习提供了更好的数据和分析的技术基础,在一些应用领域,机器已经达到甚至超过人类的智能水平,从而引发机器学习在金融、智能制造、零售、电子商务、电信等众多行业的广泛应用。

大数据应用和人工智能发展已经引起了全球关注,企业对机器学习的人才需求增大,与之密切相关的数据科学家、数据挖掘工程师、大数据分析师、机器学习分析师等数据分析类人才成为 21 世纪最有吸引力的人才。根据相关企业估计,上述人才的需求增长将持续 6~8 年。为此,教育部批准了一批高校成立数据科学与大数据技术、大数据应用、智能科学与技术等相关专业,为业界培养相应的专业人才。

机器学习是一门理论和实践并重的课程,内容比较多,很多算法也有一定的难度。此外,机器学习的应用需要一定的经验和技巧。编写一本兼顾机器学习理论的系统性,又能体现机器学习的应用,使之适合新工科教学的机器学习教材,就成为本书编写的目标。

目前,市场上已有多种机器学习的专业图书和教材,大多数图书有两种倾向,要么理论内容繁杂,对初学者和任课教师来说挑战较大,短短的几十学时难以消化;要么偏重应用和实践,理论的系统性不够,理论与应用方面的内容不够平衡。针对上述问题,作者参阅了大量文献资料,结合过去多年数据分析的研究和实践,重新梳理了机器学习的整个课程体系,使内容基本覆盖机器学习的基础内容,深入浅出,读者在此基础上可以钻研机器学习的高级算法。

本书具有以下特点。

(1) 大多数章节都有典型的 Python 算法和案例，深入浅出地解释理论，方便学习理解。本书最后附有主要参考文献，方便读者加深对教材内容的思考。此外，每章还配有思考题，以此检验读者对基本知识的理解和应用能力。

(2) 在介绍传统的机器学习理论的基础上，突出了机器学习目前主流的一些内容，包括深度学习的典型算法与应用、知识图谱、机器学习在电子推荐技术的应用等。

(3) 本书第 14 章的华为 FusionInsight 平台提供 3 个月的免费使用期，请读者根据实际情况开通使用。

本书在写作过程中，得到了教育部高等学校计算机类专业教学指导委员会-华为 ICT 产学合作项目组的大力支持，另外，研究生于召鑫、蒲实、朱荣斌、耿甲、袁雪如、陈子轩等在资料收集过程中做了一些工作，在此一并表示感谢。

赵卫东

2018 年 2 月

复旦大学

# 目 录 CONTENTS

## 第 1 章 机器学习概述 ..... 1

- 1.1 机器学习简介 ..... 2
  - 1.1.1 机器学习简史 ..... 2
  - 1.1.2 机器学习主要流派 ..... 3
- 1.2 机器学习、人工智能和数据挖掘 ..... 5
  - 1.2.1 什么是人工智能 ..... 5
  - 1.2.2 什么是数据挖掘 ..... 6
  - 1.2.3 机器学习、人工智能与数据挖掘的关系 ..... 6
- 1.3 典型机器学习应用领域 ..... 7
- 1.4 机器学习算法 ..... 13
- 1.5 机器学习的一般流程 ..... 20
- 习题 ..... 21

## 第 2 章 机器学习基本方法 ..... 23

- 2.1 统计分析 ..... 24
  - 2.1.1 统计基础 ..... 24
  - 2.1.2 常见概率分布 ..... 29
  - 2.1.3 参数估计 ..... 30
  - 2.1.4 假设检验 ..... 32
  - 2.1.5 线性回归 ..... 33
  - 2.1.6 逻辑回归 ..... 35
  - 2.1.7 判别分析 ..... 37
  - 2.1.8 非线性模型 ..... 38
- 2.2 高维数据降维 ..... 39
  - 2.2.1 主成分分析 ..... 39
  - 2.2.2 奇异值分解 ..... 42
  - 2.2.3 线性判别分析 ..... 43
  - 2.2.4 局部线性嵌入 ..... 46
  - 2.2.5 拉普拉斯特征映射 ..... 47
- 2.3 特征工程 ..... 49

- 2.3.1 特征构建 ..... 49
- 2.3.2 特征选择 ..... 50
- 2.3.3 特征提取 ..... 51
- 2.4 模型训练 ..... 51
  - 2.4.1 模型训练常见术语 ..... 51
  - 2.4.2 训练数据收集 ..... 51
- 2.5 可视化分析 ..... 52
  - 2.5.1 可视化分析的作用 ..... 53
  - 2.5.2 可视化分析方法 ..... 53
  - 2.5.3 可视化分析常用工具 ..... 54
  - 2.5.4 常见的可视化图表 ..... 56
  - 2.5.5 可视化分析面临的挑战 ..... 66
- 习题 ..... 66

## 第 3 章 决策树与分类算法 ..... 68

- 3.1 决策树算法 ..... 69
  - 3.1.1 分支处理 ..... 70
  - 3.1.2 连续属性离散化 ..... 76
  - 3.1.3 过拟合问题 ..... 78
  - 3.1.4 分类效果评价 ..... 83
- 3.2 集成学习 ..... 87
  - 3.2.1 装袋法 ..... 87
  - 3.2.2 提升法 ..... 88
  - 3.2.3 GBDT ..... 90
  - 3.2.4 随机森林 ..... 91
- 3.3 决策树应用 ..... 93
- 习题 ..... 96

## 第 4 章 聚类分析 ..... 97

- 4.1 聚类分析概念 ..... 98
  - 4.1.1 聚类方法分类 ..... 98
  - 4.1.2 良好聚类算法的特征 ..... 99

4.2 聚类分析的度量	100	5.4.1 文本分词	151
4.2.1 外部指标	100	5.4.2 命名实体识别	154
4.2.2 内部指标	101	5.4.3 词义消歧	155
4.3 基于划分的聚类	103	5.5 句法分析	155
4.3.1 $k$ -均值算法	103	5.6 语义分析	157
4.3.2 $k$ -medoids 算法	108	5.7 文本分析应用	158
4.3.3 $k$ -prototype 算法	108	5.7.1 文本分类	159
4.4 基于密度的聚类	109	5.7.2 信息抽取	161
4.4.1 DBSCAN 算法	109	5.7.3 问答系统	162
4.4.2 OPTICS 算法	111	5.7.4 情感分析	163
4.4.3 DENCLUE 算法	112	5.7.5 自动摘要	164
4.5 基于层次的聚类	115	习题	165
4.5.1 BIRCH 聚类	115	<b>第 6 章 神经网络</b>	<b>166</b>
4.5.2 CURE 算法	118	6.1 神经网络介绍	167
4.6 基于网格的聚类	121	6.1.1 前馈神经网络	167
4.7 基于模型的聚类	121	6.1.2 反馈神经网络	169
4.7.1 概率模型聚类	121	6.1.3 自组织神经网络	172
4.7.2 模糊聚类	126	6.2 神经网络相关概念	173
4.7.3 Kohonen 神经网络聚类	126	6.2.1 激活函数	173
习题	132	6.2.2 损失函数	176
<b>第 5 章 文本分析</b>	<b>134</b>	6.2.3 学习率	178
5.1 文本分析介绍	135	6.2.4 过拟合	180
5.2 文本特征提取及表示	135	6.2.5 模型训练中的问题	181
5.2.1 TF-IDF	136	6.2.6 神经网络效果评价	184
5.2.2 信息增益	136	6.3 神经网络应用	184
5.2.3 互信息	137	习题	188
5.2.4 卡方统计量	138	<b>第 7 章 贝叶斯网络</b>	<b>189</b>
5.2.5 词嵌入	138	7.1 贝叶斯理论概述	190
5.2.6 语言模型	139	7.2 贝叶斯概率基础	190
5.2.7 向量空间模型	141	7.2.1 概率论	190
5.3 知识图谱	142	7.2.2 贝叶斯概率	191
5.3.1 知识图谱相关概念	143	7.3 朴素贝叶斯分类模型	192
5.3.2 知识图谱的存储	144	7.4 贝叶斯网络推理	195
5.3.3 知识图谱挖掘与计算	145	7.5 贝叶斯网络的应用	200
5.3.4 知识图谱的构建过程	146	7.5.1 中文分词	200
5.4 词法分析	151	7.5.2 机器翻译	201

7.5.3 故障诊断	201	11.3 深度学习流行框架	264
7.5.4 疾病诊断	202	习题	265
习题	204	<b>第 12 章 高级深度学习</b>	<b>266</b>
<b>第 8 章 支持向量机</b>	<b>205</b>	12.1 高级卷积神经网络	267
8.1 支持向量机模型	206	12.1.1 目标检测与追踪	267
8.1.1 核函数	206	12.1.2 目标分割	270
8.1.2 模型原理分析	207	12.2 高级循环神经网络应用	272
8.2 支持向量机应用	210	12.2.1 Encoder-Decoder 模型	272
习题	215	12.2.2 注意力模型	273
<b>第 9 章 进化计算</b>	<b>216</b>	12.2.3 LSTM 高级应用	274
9.1 遗传算法的基础	217	12.3 无监督式深度学习	275
9.1.1 基因重组与基因突变	217	12.3.1 深度信念网络	275
9.1.2 遗传算法实现技术	218	12.3.2 生成对抗网络模型	277
9.1.3 遗传算法应用案例	222	12.4 强化学习	277
9.2 蚁群算法	223	12.5 迁移学习	279
9.3 蜂群算法	225	12.6 对偶学习	282
习题	227	习题	283
<b>第 10 章 分布式机器学习</b>	<b>229</b>	<b>第 13 章 推荐系统</b>	<b>284</b>
10.1 分布式机器学习基础	230	13.1 推荐系统基础	285
10.1.1 参数服务器	230	13.1.1 推荐系统的应用场景	285
10.1.2 分布式并行计算类型	231	13.1.2 相似度计算	286
10.2 分布式机器学习框架	232	13.2 推荐系统通用模型	288
10.3 并行决策树	238	13.2.1 推荐系统结构	288
10.4 并行 $k$ -均值算法	238	13.2.2 基于人口统计学的推荐	288
习题	240	13.2.3 基于内容的推荐	289
<b>第 11 章 深度学习</b>	<b>242</b>	13.2.4 基于协同过滤的推荐算法	290
11.1 卷积神经网络	243	13.2.5 基于图的模型	292
11.1.1 卷积神经网络简介	243	13.2.6 基于关联规则的推荐	293
11.1.2 卷积神经网络的结构	244	13.2.7 基于知识的推荐	299
11.1.3 常见卷积神经网络	246	13.2.8 基于标签的推荐	300
11.2 循环神经网络	254	13.3 推荐系统评测	301
11.2.1 RNN 基本原理	254	13.3.1 评测方法	301
11.2.2 长短期记忆网络	260	13.3.2 评测指标	302
11.2.3 门限循环单元	263	13.4 推荐系统常见问题	306
		13.5 推荐系统实例	309
		习题	318



**第 14 章 实验 ..... 319**

- 14.1 华为 FusionInsight 产品平台介绍 ..... 320
- 14.2 银行定期存款业务预测 ..... 321
  - 14.2.1 上传银行客户及存贷款数据 ..... 322
  - 14.2.2 准备存款业务分析工作区 ..... 322
  - 14.2.3 创建数据挖掘流程 ..... 323
  - 14.2.4 定期存款业务模型保存和应用 ..... 330
- 14.3 客户分群 ..... 333

- 14.3.1 分析业务需求 ..... 333
- 14.3.2 上传客户信息数据 ..... 335
- 14.3.3 准备客户分群工作区 ..... 336
- 14.3.4 创建数据挖掘流程 ..... 337
- 14.3.5 客户分群模型保存和应用 ..... 344

**附录 《机器学习》配套实验课程  
方案简介 ..... 347**

**参考文献 ..... 348**

# 01

## 第1章 机器学习概述

随着大数据的发展和计算机运算能力的不断提升,人工智能在最近几年取得了令人瞩目的成就。目前在很多行业中,都有企业开始应用机器学习技术,从而获取更深刻的洞察,为企业经营或日常生活提供帮助,提升产品服务水平。机器学习已经广泛应用于数据挖掘、搜索引擎、电子商务、自动驾驶、图像识别、量化投资、自然语言处理、计算机视觉、医学诊断、信用卡欺诈检测、证券金融市场分析、游戏和机器人等领域,机器学习相关技术的进步促进了人工智能在各个领域的发展。