

今天不研究机器学习，明天就被机器取代
你还在等什么？

Python+TensorFlow 机器学习实战

李鸥 编著

- 很系统：讲解19种机器学习经典算法，依次击破重难点
- 很图示：书中包括113张图解说明，方便读者理解
- 很实用：囊括文本识别、语音识别、图形识别、人脸识别等
- 很实战：31个实例、13个案例，详解TensorFlow机器学习

清华大学出版社

Python+TensorFlow

机器学习实战

李鸥 编著



清华大学出版社

北·京

内 容 简 介

本书通过开发实例和项目案例,详细介绍 TensorFlow 开发所涉及的主要内容。书中的每个知识点都通过实例进行通俗易懂的讲解,便于读者轻松掌握有关 TensorFlow 开发的内容和技巧,并能够得心应手地使用 TensorFlow 进行开发。

本书内容共分为 11 章,首先介绍 TensorFlow 的基本知识,通过实例逐步深入地讲解线性回归、支持向量机、神经网络算法和无监督学习等常见的机器学习算法模型。然后通过 TensorFlow 在自然语言文本处理、语音识别、图形识别和人脸识别等方面的成功应用讲解 TensorFlow 的实际开发过程。

本书适合有一定 Python 基础的工程师阅读;对于有一定基础的读者,可通过本书快速地将 TensorFlow 应用到实际开发中;对于高等院校的学生和培训机构的学员,本书也是入门和实践机器学习的优秀教材。

本书对应的电子课件和实例源代码可以到 <http://www.tupwk.com.cn/downpage> 下载,也可通过扫描前言中的二维码下载。

本书封面贴有清华大学出版社防伪标签,无标签者不得销售。

版权所有,侵权必究。侵权举报电话:010-62782989 13701121933

图书在版编目(CIP)数据

Python+TensorFlow机器学习实战 / 李鸥 编著. —北京:清华大学出版社, 2019
ISBN 978-7-302-52260-7

I. ①P… II. ①李… III. ①软件工具—程序设计②人工智能—算法 IV. ①TP311.561②TP18

中国版本图书馆 CIP 数据核字(2019)第 018842 号

责任编辑:胡辰浩

装帧设计:孔祥峰

责任校对:牛艳敏

责任印制:沈 露

出版发行:清华大学出版社

网 址: <http://www.tup.com.cn>, <http://www.wqbook.com>

地 址:北京清华大学学研大厦 A 座 邮 编:100084

社总机:010-62770175 邮 购:010-62786544

投稿与读者服务:010-62776969, c-service@tup.tsinghua.edu.cn

质量反馈:010-62772015, zhiliang@tup.tsinghua.edu.cn

印装者:清华大学印刷厂

经 销:全国新华书店

开 本:185mm×260mm 印 张:15.5 字 数:358 千字

版 次:2019 年 6 月第 1 版 印 次:2019 年 6 月第 1 次印刷

印 数:1~3000

定 价:79.00 元

产品编号:080563-01

前言

2016年3月，谷歌公司的AlphaGo与职业九段棋手李世石进行了围棋人机大战，最终AlphaGo以4比1的总比分获胜，这引起了全球对人工智能的热议。同时，百度推出的无人驾驶，科大讯飞推出的“语音识别”，以及高铁进站的人脸识别的广泛应用，将机器学习转变为信息科技企业的研究与应用的常见内容，这也让我们的日常生活更为便捷。

其实，机器学习已经走过符号主义时代、概率论时代、联结主义时代，从最初的仅是专家研究的数学理论、经典算法，逐步发展并蜕变为可以为大部分项目直接使用的平台框架。

2015年11月9日，谷歌在GitHub上开源了TensorFlow框架，该框架是谷歌的机器学习框架，具有高度的灵活性和可移植性。在TensorFlow中，将各种经典算法特别是神经网络模型组织成一个平台，能够让我们更便捷地在目标领域实践机器学习算法。

TensorFlow作为最流行的机器学习框架之一，具有对Python语言的良好支持，这有效降低了进行机器学习开发的门槛，让更多的工程师能够以低成本投身到人工智能的浪潮中。TensorFlow框架能够支持CPU、GPU或Google TPU等硬件环境，让机器学习能够便捷地移植到各种环境中。

本书将全面阐述TensorFlow机器学习框架的原理、概念，详细讲解线性回归、支持向量机、神经网络算法和无监督学习等常见的机器学习算法模型，并通过TensorFlow在自然语言文本处理、语音识别、图形识别和人脸识别等方面的成功应用来讲解TensorFlow的实际开发过程。本书在语言上力求幽默直白、轻松活泼，避免云山雾罩、晦涩难懂。在讲解形式上图文并茂，由浅入深，抽丝剥茧。通过阅读本书，读者可以少走很多弯路，快速上手TensorFlow开发。

本书特色

1. 内容丰富、全面

全书内容共分11章，从机器学习概述到TensorFlow基础，再到实际应用，内容几乎涵盖TensorFlow开发的所有方面。

2. 实例丰富、案例典型、实用性强

本书对每一个知识点都以实际应用的形式进行讲解，帮助读者理解和掌握相关的开发技术。本书还在最后提供了TensorFlow在图形识别、文本识别和语音识别等方面成功应用的实例，帮助读者提高实战水平。

3. 紧跟技术趋势

本书针对目前发布的TensorFlow的常用版本1.3进行讲解，并涉及1.6版本的变化，摒弃了以前版本中不再使用的功能，以适应技术的发展趋势。

4. 举一反三

本书写作由浅入深、从易到难，并注意知识点之间的联系，让读者掌握一个知识点后，能够触类旁通、举一反三，编写相应的代码。

本书内容及体系结构

第1章简单讲述机器学习的发展、分类以及经典算法，介绍TensorFlow的发展和优势，并详细介绍不同操作系统环境下TensorFlow开发环境的准备过程。

第2章讲解TensorFlow的基础知识，包括基础框架、源代码结构、基础概念，并通过运行一个官方示例展示了可视化的调试。

第3章讲解TensorFlow在实际进行机器学习时的加载训练数据、构建训练模型、进行数据训练、评估和预测四大步骤中常用的方法和技巧。

第4章详细讲解机器学习算法中最基础的线性模型：回归模型和逻辑回归模型。

第5章讲解TensorFlow中支持向量机算法的基本原理及核函数，并使用SVM完成线性回归拟合、逻辑回归分类以及非线性数据分类等。

第6章对神经网络模型进行详细介绍，讲解神经元模型、神经网络层等基本原理解，并讲解全连接神经网络、卷积神经网络和循环神经网络等主要神经网络的原理与计算过程，并在TensorFlow中使用具体案例讲解通用神经网络层的构建、卷积层的使用、池化层的使用、循环神经元的构建以及损失函数的选择等。

第7章主要介绍无监督学习的概念和经典算法。

第8章讲解TensorFlow在自然语言文本处理中的应用，如学写唐诗、影评分类以及智能聊天机器人等。

第9章讲解TensorFlow在语音处理方面的应用，如听懂数字、听懂中文以及语音合成等。

第10章讲解TensorFlow在图像处理方面的应用，如图像处理中的物体识别与检测、图像描述。

第11章讲解TensorFlow在人脸识别方面的应用，介绍人脸识别的原理和分类、人脸比对以及从人脸判别性别和年龄。

本书读者对象

- 初中级程序员。
- 高等院校师生。
- 培训机构学员。
- 希望使用机器学习的工程师。

致谢

在本书的成稿过程中，熊诺亚对书稿的完整性和系统性提出了宝贵的意见，在此，特别表示感谢。

本书对应的电子课件和实例源代码可以到<http://www.tupwk.com.cn/downpage>下载，也可通过扫描下方的二维码下载。



编著者

目录

第1章 机器学习概述

- 1.1 人工智能 1
- 1.2 机器学习 2
 - 1.2.1 机器学习的发展 2
 - 1.2.2 机器学习的分类 3
 - 1.2.3 机器学习的经典算法 4
 - 1.2.4 机器学习入门 6
- 1.3 TensorFlow简介 6
 - 1.3.1 主流框架的对比 7
 - 1.3.2 TensorFlow的发展 9
 - 1.3.3 使用TensorFlow的公司 10
- 1.4 TensorFlow环境准备 10
 - 1.4.1 Windows环境 11
 - 1.4.2 Linux环境 21
 - 1.4.3 Mac OS环境 22
- 1.5 常用的第三方模块 22
- 1.6 本章小结 23

第2章 TensorFlow基础

- 2.1 TensorFlow基础框架 24
 - 2.1.1 系统框架 24
 - 2.1.2 系统的特性 26
 - 2.1.3 编程模型 27
 - 2.1.4 编程特点 28
- 2.2 TensorFlow源代码结构分析 30
 - 2.2.1 源代码下载 30
 - 2.2.2 TensorFlow目录结构 30

- 2.2.3 重点目录 31
- 2.3 TensorFlow基本概念 33
 - 2.3.1 Tensor 33
 - 2.3.2 Variable 34
 - 2.3.3 Placeholder 35
 - 2.3.4 Session 36
 - 2.3.5 Operation 36
 - 2.3.6 Queue 37
 - 2.3.7 QueueRunner 38
 - 2.3.8 Coordinator 39
- 2.4 第一个TensorFlow示例 40
 - 2.4.1 典型应用 41
 - 2.4.2 运行TensorFlow示例 43
- 2.5 TensorBoard可视化 45
 - 2.5.1 SCALARS面板 45
 - 2.5.2 GRAPHS面板 47
 - 2.5.3 IMAGES面板 48
 - 2.5.4 AUDIO面板 49
 - 2.5.5 DISTRIBUTIONS面板 49
 - 2.5.6 HISTOGRAMS面板 49
 - 2.5.7 PROJECTOR面板 50
- 2.6 本章小结 50

第3章 TensorFlow进阶

- 3.1 加载数据 51
 - 3.1.1 预加载数据 51
 - 3.1.2 填充数据 51

3.1.3	从CSV文件读取数据	52	5.2.3	进行数据训练	81
3.1.4	读取TFRecords数据	54	5.2.4	运行总结	82
3.2	存储和加载模型	58	5.3	拟合逻辑回归	83
3.2.1	存储模型	58	5.3.1	生成训练数据	83
3.2.2	加载模型	59	5.3.2	定义训练模型	84
3.3	评估和优化模型	60	5.3.3	进行数据训练	85
3.3.1	评估指标的介绍与使用	60	5.3.4	运行总结	86
3.3.2	模型调优的主要方法	61	5.4	非线性二值分类	87
3.4	本章小结	63	5.4.1	生成训练数据	87
第4章 线性模型			5.4.2	定义训练模型	88
4.1	常见的线性模型	64	5.4.3	进行数据训练	89
4.2	一元线性回归	65	5.4.4	运行总结	89
4.2.1	生成训练数据	65	5.5	非线性多类分类	91
4.2.2	定义训练模型	66	5.5.1	生成训练数据	91
4.2.3	进行数据训练	66	5.5.2	定义训练模型	92
4.2.4	运行总结	67	5.5.3	进行数据训练	93
4.3	多元线性回归	68	5.5.4	运行总结	94
4.3.1	二元线性回归算法简介	68	5.6	本章小结	95
4.3.2	生成训练数据	69	第6章 神经网络		
4.3.3	定义训练模型	70	6.1	神经网络简介	96
4.3.4	进行数据训练	70	6.1.1	神经元模型	97
4.3.5	运行总结	70	6.1.2	神经网络层	100
4.4	逻辑回归	71	6.2	拟合线性回归问题	102
4.4.1	逻辑回归算法简介	71	6.2.1	生成训练数据	102
4.4.2	生成训练数据	73	6.2.2	定义神经网络模型	102
4.4.3	定义训练模型	74	6.2.3	进行数据训练	103
4.4.4	进行数据训练	74	6.2.4	运行总结	104
4.4.5	运行总结	75	6.3	MNIST数据集	104
4.5	本章小结	76	6.3.1	MNIST数据集简介	105
第5章 支持向量机			6.3.2	数据集图片文件	105
5.1	支持向量机简介	77	6.3.3	数据集标记文件	106
5.1.1	SVM基本型	77	6.4	全连接神经网络	106
5.1.2	SVM核函数简介	79	6.4.1	加载MNIST训练数据	106
5.2	拟合线性回归	80	6.4.2	构建神经网络模型	107
5.2.1	生成训练数据	80	6.4.3	进行数据训练	108
5.2.2	定义训练模型	81	6.4.4	评估模型	109

6.4.5	构建多层神经网络模型	110	第8章	自然语言文本处理	
6.4.6	可视化多层神经网络模型	111	8.1	自然语言文本处理简介	152
6.5	卷积神经网络	113	8.1.1	处理模型的选择	152
6.5.1	卷积神经网络简介	114	8.1.2	文本映射	153
6.5.2	卷积层	115	8.1.3	TensorFlow文本处理的 一般步骤	156
6.5.3	池化层	119	8.2	学写唐诗	157
6.5.4	全连接神经网络层	121	8.2.1	数据预处理	157
6.5.5	卷积神经网络的发展	121	8.2.2	生成训练模型	158
6.6	通过卷积神经网络处理MNIST	122	8.2.3	评估模型	160
6.6.1	加载MNIST训练数据	122	8.3	智能影评分类	163
6.6.2	构建卷积神经网络模型	123	8.3.1	CBOW嵌套模型	163
6.6.3	进行数据训练	127	8.3.2	构建影评分类模型	167
6.6.4	评估模型	127	8.3.3	训练评估影评分类模型	169
6.7	循环神经网络	128	8.4	智能聊天机器人	170
6.7.1	循环神经网络简介	128	8.4.1	Attention机制的Seq2Seq 模型	170
6.7.2	基本循环神经网络	129	8.4.2	数据预处理	173
6.7.3	长短期记忆网络	131	8.4.3	构建智能聊天机器人模型	174
6.7.4	双向循环神经网络简介	134	8.4.4	训练模型	177
6.8	通过循环神经网络处理MNIST	135	8.4.5	评估模型	179
6.8.1	加载MNIST训练数据	136	8.5	本章小结	180
6.8.2	构建神经网络模型	136	第9章	语音处理	
6.8.3	进行数据训练及评估模型	137	9.1	语音处理简介	181
6.9	递归神经网络	138	9.1.1	语音识别模型	181
6.9.1	递归神经网络简介	138	9.1.2	语音合成模型	183
6.9.2	递归神经网络的应用	139	9.2	听懂数字	183
6.10	本章小结	140	9.2.1	数据预处理	184
第7章	无监督学习		9.2.2	构建识别模型	185
7.1	无监督学习简介	141	9.2.3	训练模型	185
7.1.1	聚类模型	141	9.2.4	评估模型	185
7.1.2	自编码网络模型	142	9.3	听懂中文	185
7.2	K均值聚类	142	9.3.1	数据预处理	186
7.2.1	K均值聚类算法简介	142	9.3.2	构建识别模型	188
7.2.2	K均值聚类算法实践	144	9.3.3	训练模型	191
7.3	自编码网络	147	9.3.4	评估模型	191
7.3.1	自编码网络简介	147	9.4	语音合成	192
7.3.2	自编码网络实践	148			
7.4	本章小结	151			

9.4.1 Tacotron模型	192	10.5.1 看图说话原理	218
9.4.2 编码器模块	193	10.5.2 看图说话模型的构建	218
9.4.3 解码器模块	196	10.5.3 看图说话模型的训练	220
9.4.4 后处理模块	197	10.5.4 评估模型	221
9.5 本章小结	197	10.6 本章小结	222
第10章 图像处理		第11章 人脸识别	
10.1 机器学习的图像处理简介	198	11.1 人脸识别简介	223
10.1.1 图像修复	198	11.1.1 人脸图像采集	223
10.1.2 图像物体识别与检测	199	11.1.2 人脸检测	224
10.1.3 图像问答	201	11.1.3 人脸图像预处理	224
10.2 图像物体识别	201	11.1.4 人脸关键点检测	224
10.2.1 数据预处理	201	11.1.5 人脸特征提取	224
10.2.2 生成训练模型	203	11.1.6 人脸比对	225
10.2.3 训练模型	205	11.1.7 人脸属性检测	225
10.2.4 评估模型	206	11.2 人脸验证	225
10.3 图片验证码识别	208	11.2.1 数据预处理	226
10.3.1 验证码的生成	208	11.2.2 运行FaceNet模型	226
10.3.2 数据预处理	209	11.2.3 实现人脸验证	229
10.3.3 生成训练模型	211	11.3 性别和年龄的识别	231
10.3.4 训练模型	212	11.3.1 Adience数据集	231
10.3.5 评估模型	213	11.3.2 数据预处理	232
10.4 图像物体检测	214	11.3.3 生成训练模型	233
10.4.1 物体检测系统	214	11.3.4 训练模型	235
10.4.2 物体检测系统实践	215	11.3.5 评估模型	236
10.5 看图说话	217	11.4 本章小结	237

第1章

机器学习概述

本章介绍人工智能和机器学习的发展，讲解机器学习的主要框架，解释TensorFlow的作用、特性以及开发环境的准备过程。

1.1 人工智能

毫无疑问，目前人工智能在全球的火热与AlphaGo(阿尔法狗)的战绩密不可分。2016年3月，谷歌公司的AlphaGo与职业九段棋手李世石进行围棋人机大战，最终以4比1的总比分获胜，这引起了全球热议。2017年年初，AlphaGo化身为Master，在棋类平台上横扫中日韩围棋高手，取得60连胜，再度引发全民对人工智能的讨论。

虽然人工智能是在AlphaGo战胜李世石之后才成了坊间谈资，引起所有人的关注，但人工智能的提出已经有近百年的历史。

早在20世纪50年代，计算机科学家就提出了“人工智能”的概念，想制造出和人类外形相同、能够与人类正常对话、能够自我学习的机器。阿兰·图灵还提出了著名的“图灵测试”来判定计算机是否智能：如果一台机器能够与人类展开对话而不被辨别出其机器身份，那么称这台机器具有智能。从此以后，人工智能就一直是人们在科研、工业以及电影中努力实现的目标，也确实在不断地发展。

现在，人工智能已经发展为一门广泛的交叉和前沿科学，涉及计算机科学、心理学、哲学和语言学等学科，也被广泛地应用到语音识别、图像识别、自然语言处理等领域。

在国际上，谷歌、微软、IBM等都有自己的人工智能项目，如谷歌的DeepMind、IBM的Watson、微软的Torque等项目。

国内的各大公司也积极投身于人工智能领域。百度成立了Apollo基金和DuerOS基金，推动中国AI的发展；腾讯创建了人工智能实验室AI Lab，专注于人工智能的基础研究；阿里巴巴成立的人工智能实验室，主要面向消费级的AI产品研发；搜狗向清华大学捐赠1.8

亿元，一起成立了“天工智能计算研究院”等。目前这些公司也陆续推出了各自的产品：腾讯开发的机器人“Dreamwriter”，百度的无人驾驶，搜狗、科大讯飞等公司的“语音识别”，旷视科技的“Face++”人脸识别等。

1.2 机器学习



1.2.1 机器学习的发展

为了让计算机能够实现类似人类的智能，在计算机的实际实现上出现了两种完全不同的方向：一种是采用传统的编程技术，使系统呈现智能的效果；另一种是采用计算机训练学习的方式来实现智能的效果。一般来说，现在我们使用的机器学习都是通过算法来解析数据、学习数据的，然后据此对真实世界中的事件做出决策和预测。

“机器学习”这一术语由IBM的科学家亚瑟·塞缪尔提出。他在1952年开发了一个跳棋程序，该程序能够观察当前位置，并学习一个隐含的模型，从而为后续动作提供更好的指导，并且随着该程序运行时间的增加，可以实现越来越可靠的后续指导。他针对这种计算机的实现能力提出了“机器学习”。

在机器学习领域，计算机科学家不断探索，基于不同的理论创建出不同的机器学习模型。从发展历程来说，大致经历了三个阶段：符号主义时代、概率论时代以及联结主义时代。

- 符号主义时代(1980年左右)。以知识工程为主要理论依据，使用服务器或大型机进行架构运算，通过符号、规则和逻辑来表征知识和进行逻辑推理，常用的算法有规则和决策树，实用性有限。
- 概率论时代(1990—2000年)。以概率论为主要理论依据，使用小型服务器集群进行架构运算，通过获取发生的可能性来进行概率推理，常用的算法有朴素贝叶斯或马尔可夫算法，具有可扩展的比较或对比功能，对许多任务都表现得足够好。
- 联结主义时代(2010年左右)。以神经科学和概率为主要理论依据，使用云计算架构，通过使用概率矩阵和加权神经元来动态地识别和归纳模式，常用的算法有神经网络，能够让计算机“看懂图像”“听懂语言”，甚至能够分析人类在语言背后表达的情绪。

不同算法在不同应用场景下有着不同的表现，每一个阶段仅仅取得了某些领域的突破性进展，并没有完全颠覆前一阶段的成果。相信在后续的发展中，将会把符号规则理论、概率论、神经科学和进化论等理论相融合，并演变出不同的算法，通过多种学习方式获得知识或经验，推动机器学习继续发展。

1.2.2 机器学习的分类

机器学习是计算机进行数据处理,找到数据间映射关系的过程。在进行数据处理分析时,对于输入的数据,有的是经过人工来定义数据标签的,方法是先找到数据的特征与其标签的映射关系,再凭借这种映射关系,对未进行标签定义的数据进行标签定义。输入的初始数据也有一些是没有经过人工定义数据标签的,只是单纯依靠数据处理分析来找到数据之间的标签映射关系。

可以按照输入的数据本身是否已被标定特定的标签将机器学习区分为有监督学习、无监督学习以及半监督学习三类。

1. 有监督学习

有监督学习(Supervised Learning)就是样本数据集中的数据,包括样本数据以及样本数据的标签。

进行学习的目的就是找到样本数据与样本数据标签的映射关系。通过对样本数据的不断学习、不断修正学习中的偏差,使得找到的映射关系更准确,从而不断提高学习的准确率。当学习完成后,再给予新的未知数据,能够依据学习的映射关系计算出相对正确的结果。由于样本数据中既包括数据也包括标签,因此训练的效果往往都不错。

有监督学习主要用于解决两大类问题:回归问题(Regression Problem)和分类问题(Classification Problem)。

回归问题就是通过对现有数据的分析,找到映射关系,对以后的事情进行预测的情况。比如,我们想预测未来房价会是多少。我们获取以前的房价与时间的数据,可以将这些数据看作多维度坐标系中的坐标点,通过回归分析,建立数据的关系模型,求出一个最符合这些已知数据集的解析函数,然后通过这个解析函数来预估未来的房价。对于解决回归问题,主要有线性回归(Linear Regression)、决策树(Decision Tree)、随机森林(Random Forest)、梯度提升决策树(Gradient Boosting Tree)、神经网络(Neural Network)等算法可供使用。

分类问题就是通过对现有数据的分析,找到数据间的联系与区别,对数据进行分类。比如,判断某地房价的“涨”与“跌”的问题。我们获取以前的房价、地区、户型和时间等数据,通过这些数据建立数据与“涨”和“跌”的关系模型。当输入新的值时,能够根据关系模型判断房价是“涨”还是“跌”了。对于解决分类问题,主要有逻辑回归(Logistics Regression)、决策树(Decision Tree)、随机森林(Random Forest)、梯度提升决策树(Gradient Boosting Tree)、核函数支持向量机(Kernel SVM)、朴素贝叶斯(Naive Bayes)、SVM线性分类(Linear SVM)、神经网络(Neural Network)等算法可供使用。

2. 无监督学习

无监督学习(Unsupervised Learning)就是在样本数据中只有数据,而没有对数据进行标记。无监督学习的目的就是让计算机对这些原始数据进行分析,让计算机自己去学习、找到数据之间的某种关系。

无监督学习与有监督学习的明显区别就是在样本数据中只有数据，没有标记。由于没有对数据进行标记，因此学习的结果也难以验证是否正确，也难以对学习的模型进行正确率的判断。对于无监督学习的这种特点，学习的思路和目的主要有两类：聚类(Clustering)和强化学习(Reinforcement Learning, RL)。

聚类就是对于未标记的数据，在训练时根据数据本身的数据特征进行训练，呈现出数据集聚的形式，每一个集聚群中的数据，彼此都有相似的性质，从而形成分组。比如我们使用的今日头条，它每天会收集大量的新闻，然后把它们全部聚类，就会自动分成娱乐、科技和政治等几十个不同的组，每个组内的新闻都具有相似的内容结构。

强化学习是游戏中常用的一种学习方式，是指在学习过程中增加一种延迟奖赏机制。通过学习过程中的延迟奖赏激励函数，可以让机器学习到当前状态下，执行哪一种操作使得最终的奖赏最多，从而让机器学习获得一种类似于决策的能力，比如AlphaGo也使用了这种强化学习方式。

用于无监督学习的经典算法有聚类算法、EM算法和深度学习算法等。

3. 半监督学习

半监督学习(Semi-Supervised Learning)是介于有监督学习和无监督学习之间的学习。一般来说，在半监督学习输入的数据样本中，存在一部分进行了标记的数据，但是大量存在的是没有进行标记的数据。

为了利用未标记的样本，必须先对未标记样本揭示的数据分布信息与类别进行假设，最常见的两种假设方式是聚类假设(Cluster Assumption)和流形假设(Manifold Assumption)。

对于聚类假设，是假设数据存在簇结构，同一个簇的样本属于同一个类别。对标记数据和未标记数据进行聚类，如果待预测样本与标记样本聚在一起，则认为待预测样本属于标记样本类。

对于流形假设，是假设数据分布在一种流形结构上，邻近的样本拥有相似的输出值。邻近的程度常用相似程度进行刻画。流形假设对输出值没有限制，相对于聚类假设而言，它的适用范围更广，可用于更多类型的学习任务。

1.2.3 机器学习的经典算法

随着机器学习的不断发展，出现了许多经典算法，这些算法为我们解决实际问题提供了强大的支持。

1. 线性模型

线性模型就是使用简单的公式通过一组数据点来查找最优拟合线。然后，通过已知的变量方程，求出需要预测的变量。对于不同形式的线性模型算法，主要包括线性回归(Linear Regression)和逻辑回归(Logistic Regression)。

线性回归从二维几何平面的角度可以理解为，在平面中存在已知的数据点，通过学习处理，找到一条线能够建立这些点之间的关系的模型。线性回归用于解决回归问题，是最简单的线性模型，易于理解。同时，由于模型太简单而不能反映变量之间复杂的关系，因

此容易出现过拟合的情形。

逻辑回归是给定样本属于类别“1”和类别“-1”的概率，用于解决分类问题，与线性回归的特点一样，易于理解但无法反映变量间的复杂关系，易出现过拟合的情形。

2. 树型模型

树型模型用于探索数据集中数据的特性，并且能够对数据按照数据特征进行分类处理，可以用于解决分类和回归问题。树型模型高度精确、稳定且易于解释，可以映射非线性关系以求解问题，主要包括决策树(Decision Tree)、随机森林(Random Forest)和梯度提升决策树(Gradient Boosting Tree)。

决策树是使用分支方法来显示决策的每个可能结果的图，它对所有的可能性进行梳理。这种算法易于理解和实现，但是由于决策树有时太简单，无法处理复杂的数据，因此一般不会单独使用。

随机森林是许多决策树的平均，每个决策树都用数据的随机样本训练。森林中每个独立的树都比完整的决策树弱，但是通过将它们结合在一起，可以通过多样性获得更高的整体表现。该算法非常容易构建并且表现往往良好，但是相比于其他算法输出预测可能较慢。

梯度提升决策树和随机森林类似，都是由弱决策树构成的，但最大的区别在于：梯度提升决策树中，树是一个接一个被相继训练的，每个随后的树主要用被先前树错误识别的数据进行训练。这使得梯度提升更少地集中于容易预测的情况并更多地集中于困难的情况。该算法训练速度快且表现非常好，但是训练数据即使出现小的变化，也会在模型中产生彻底的改变，因此可能会产生不可解释的结果。

3. 支持向量机

支持向量机(SVM)基于统计学理论而提出，是机器学习中一种大放光彩的经典算法。

支持向量机算法通过给予严格的优化条件获得分类界线，并且通过与高斯核等核函数的结合，通过非线性映射，把样本空间映射到高维乃至无穷维的特征空间，使得原来样本空间中非线性可分的问题转变为特征空间中线性可分的问题。它几乎不增加计算的复杂性，而且在某种程度上避免了“维数灾难”，训练较为简单，是一种广泛应用的机器学习方式。

4. 人工+神经网络

人工+神经网络算法起步较早，但是发展坎坷。

在20世纪20年代就已经提出了人工神经网络模型以及关键的反向传播算法。但是由于受当时计算机运算能力的限制，难以在多层神经网络中进行训练，通常都是只有一层隐层节点的浅层模型。这种模型的神经网络算法比较容易出现过训练现象，而且训练速度比较慢。在层次比较少的情况下，训练效果往往不如其他算法。

在2006年，Hinton提出了深度学习算法，增加了神经网络的层数和一些处理技巧。在丰富的训练数据以及强劲的计算机运行能力的帮助下，神经网络的能力大大提高。

目前，深度学习模型在目标识别、语音识别、自然语言处理等领域取得了突飞猛进的

成果，是目前最热门的机器学习方法，也是本书讲解的主要内容。但是也有使用前提，一是深度学习模型需要大量的训练数据，才能展现出神奇的效果；二是深度学习对计算能力要求更高。在有些领域采用传统的、简单的机器学习方法可以很好地解决问题，就没必要非得使用复杂的深度学习方法。

1.2.4 机器学习入门

对某个领域进行学习的第一步就是要尽快了解全貌以搭建出整体的知识体系，然后在实践中不断提升对该领域的认识。对于机器学习领域，整体的知识体系如下。

1. 数学知识

机器学习的目标是通过现有数据构建和训练模型，用于数据的分析与预测。计算机能够做的只有计算，而如何将训练过程抽象为数学函数就是需要我们掌握的能力。在现有的经典算法中涉及概率统计、矩阵运算、微积分导数等数学知识。对于这些知识学过最好，没有学过也没关系，本书会讲解在实际应用中所需要的原理和结论，其中会涉及必要的公式推导证明。

2. 编程语言

Python是一种面向对象的解释型高级编程语言，众多的机器学习框架都支持Python，因此它成了机器学习的首选语言。本书也将使用Python作为实现语言进行讲解，希望读者已经掌握了Python语言。

3. 经典机器学习理论和基本算法

经典的机器学习算法包括线性回归、逻辑回归、SVM支持向量机、神经网络算法等，以及通过各种基本算法处理数据时存在的正则化需求、过拟合现象等基本的算法特性和适用环境。本书将对这些基本算法进行详解并通过实例来说明这些算法的使用。

4. 动手实践机器学习

掌握了机器学习的基础知识后，就可以动手实践机器学习模型。首先需要选择一个开源的机器学习框架。在选择机器学习框架方面的主要考虑因素就是哪个框架的使用范围广、使用人数多。目前，TensorFlow由于由谷歌进行开源推广且有着大量的开发者群体，更新和发布速度非常快，是非常不错的选择。

1.3 TensorFlow简介

TensorFlow是谷歌公司推出的机器学习开源神器，是谷歌基于DistBelief进行研发的第二代人工智能学习系统。DistBelief是谷歌内部开发和使用的机器学习框架，但是它严重依赖于Google内部硬件，仅适用于开发神经网络算法等，因此难以广泛使用。谷歌在DistBelief的基础上，提高了运算效率、框架的灵活性和可移植性，形成了TensorFlow框

架。目前, TensorFlow已被广泛用于文本处理、语音识别和图像识别等多项机器学习和深度学习领域。

1.3.1 主流框架的对比

在机器学习的开源框架中, Google(谷歌)、Microsoft(微软)、Facebook(脸书)和Amazon(亚马逊)等巨头都有着自己的机器学习框架并进行了一定程度的开源。此外, 还有伯克利大学的贾扬清主导开发的Caffe、蒙特利尔大学Lisa Lab团队开发的Theano以及其他个人或商业组织贡献的框架。可以说, 各种开源的深度学习框架层出不穷。

1. 基本情况

TensorFlow是Google的可移植机器学习和神经网络库, 可扩展性强。TensorFlow对Python有着良好的编程语言支持, 支持CPU、GPU和Google TPU等硬件, 并且已经拥有各种各样的模型和算法, 在深度学习上有非常出色的表现。另外, TensorFlow由谷歌进行主导, 在文档和实例方面也有着良好的支持。

MXNet是亚马逊的机器学习框架, 具有较强的可移植性和可扩展性, 对Python、R、Scala、Julia和C++等编程语言有着不同程度的支持。

Deeplearning4j(DL4J)是一个专注于深度神经网络的Java库, 可以与Hadoop和其他基于Java的分布式框架集成。

Microsoft Cognitive Toolkit(CNTK)是微软的开源深度学习框架, 支持Python编程语言, 拥有各种各样的神经网络模型, 并且支持强化学习、生成对抗网络等, 是一个功能强大的工具。但是由于交流的社区小, 在文档和实例方面的学习资料很少。

Caffe深度学习项目最初是一个用于解决图像分类问题的框架, 后来逐步成长为一个强大的机器学习框架。但是由于其创始人现已离开项目, 有一段时间已不再进行更新。该项目的下一步进展不明确, 建议不再使用。

Torch是Facebook主推的机器学习框架, 基于Lua语言进行开发, 广泛支持各种机器学习模型和算法。但是由于Lua语言是机器学习中相对冷门的语言, 因此增加了学习成本。

Theano由蒙特利尔大学机器学习研究所(MILA)创建。Theano支持Python语言, 并且能够支持其他的深度学习框架。但因为它由研究机构开发, 在API方面并不完善, 若要写出效率高的Theano框架, 需要对隐藏在框架背后的算法也相当熟悉, 所以它只在研究中极为流行, 但在项目开发方面难度相对较大。

Keras由Francis Chollet编写和维护, 基于Python进行编写, 能够运行在Theano或TensorFlow上, 可以将它看成对Theano或TensorFlow的再一次封装。Keras由于水平高, 对用户友好, 因此能够更加便捷地编写卷积神经网络、递归神经网络等机器学习模型。目前, TensorFlow将Keras添加为TensorFlow核心中的高级框架, 成为TensorFlow的默认API。

对于这些主要的机器学习框架, 基本情况的对比如表1.1所示。