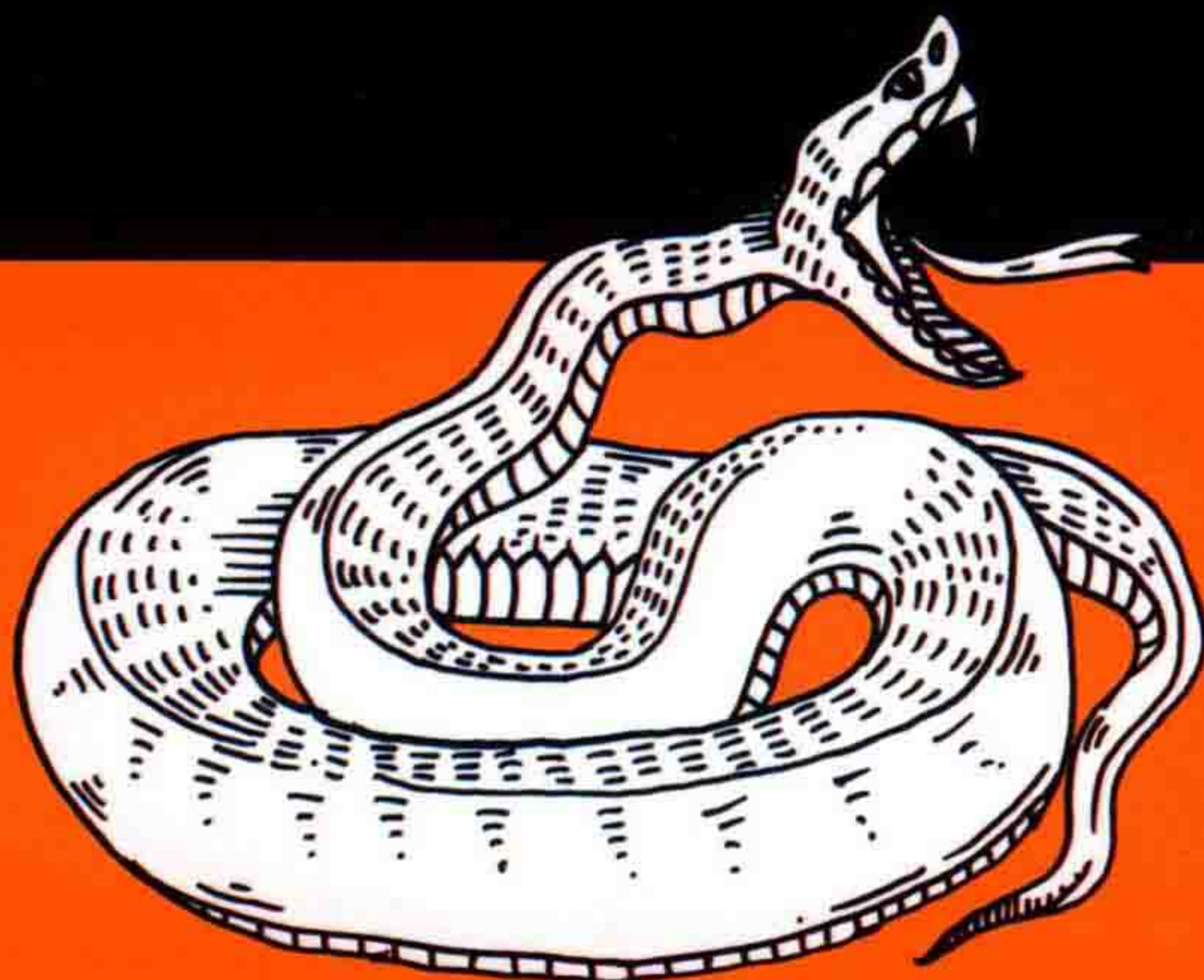


45个大型真实项目案例
95节同步微视频讲解
18小时Python 3 语法教学视频

Python 3

数据分析与 机器学习实战

龙马高新教育◎编著

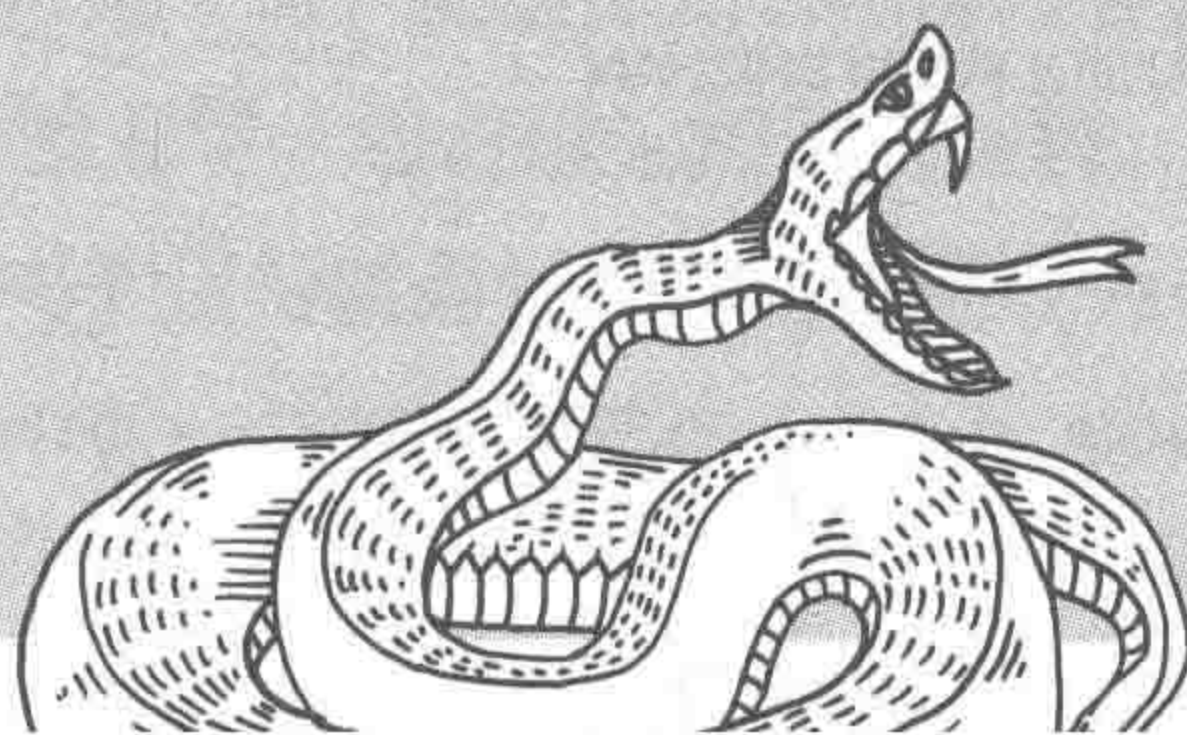


北京大学出版社
PEKING UNIVERSITY PRESS

Python 3

数据分析与 机器学习实战

龙马高新教育◎编著



北京大学出版社
PEKING UNIVERSITY PRESS

内 容 提 要

机器之所以能学习，是因为大量的数据分析。

本书首先讲述数据分析的过程，然后详细介绍常用的机器学习理论、算法与案例（大型案例 29 个），最终以解决实际问题驱动成书。本书主要介绍的机器学习算法及数据分析方法包括数据预处理、分类问题、预测问题、网络爬虫、数据降维、数据压缩、关联分析、集成学习和深度学习等。

全书共 17 章分为三大部分：第 0~3 章介绍 Python 的基础知识、安装和基本语法；第 4~7 章介绍 Python 的编程、机器学习基础和 Python 中常用的第三方库函数及数据预处理的基本方法；第 8~16 章介绍常用的机器学习分析算法及深度学习。每章都结合多个经典案例图文并茂地介绍机器学习的原理和实现方法。

本书通俗易懂，且赠送全程同步教学录像和 Python 3 编程基础录像，属学习 Python 及机器学习理论和数据分析的入门与提高课程，对于不熟悉 Python 又想学习机器学习相关算法的初学者来说非常适合。

图书在版编目 (CIP) 数据

Python 3 数据分析与机器学习实战 / 龙马高新教育编著. — 北京 : 北京大学出版社, 2018.8
ISBN 978-7-301-29566-3

I . ① P… II . ① 龙… III . ① 软件工具 — 程序设计 IV . ① TP311.561

中国版本图书馆 CIP 数据核字 (2018) 第 116684 号

- | | |
|-------|--|
| 书 名 | Python 3 数据分析与机器学习实战
PYTHON 3 SHUJU FENXI YU JIQI XUEXI SHIZHAN |
| 著作责任者 | 龙马高新教育 编著 |
| 责任编辑 | 尹毅 |
| 标准书号 | ISBN 978-7-301-29566-3 |
| 出版发行 | 北京大学出版社 |
| 地 址 | 北京市海淀区成府路 205 号 100871 |
| 网 址 | http://www.pup.cn 新浪微博 : @ 北京大学出版社 |
| 电子信箱 | pup7@pup.cn |
| 电 话 | 邮购部 62752015 发行部 62750672 编辑部 62570390 |
| 印 刷 者 | 北京大学印刷厂 |
| 经 销 者 | 新华书店 |
| | 787 毫米 × 1092 毫米 16 开本 19.75 印张 436 千字 |
| | 2018 年 8 月第 1 版 2018 年 8 月第 1 次印刷 |
| 印 数 | 1—4000 册 |
| 定 价 | 69.00 元 |

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

举报电话：010-62752024 电子信箱：fd@pup.pku.edu.cn

图书如有印装质量问题，请与出版部联系，电话：010-62756370

前言

本书专为 Python 初学者和爱好者打造,旨在使读者学会、掌握 Python 相关知识和技能,并能进行项目开发。当您认真、系统地学习本书之后,就可以骄傲地说“我是一位真正的 Python 程序员了!”,哪怕目前您还是初学者。

机器学习 (Machine Learning, ML) 是一门多领域交叉的学科,是人工智能的核心,其应用遍及人工智能的各个领域,专门研究计算机是如何模拟或实现人类的学习行为,以获取新的知识或技能,重新组织已有的知识结构,使之不断改善自身的性能。在机器学习过程中,需要使用大量数据,而数据分析是指用适当的方法对收集来的大量数据进行分析,提取有用信息和形成结论而对数据加以详细研究和概括总结的过程。

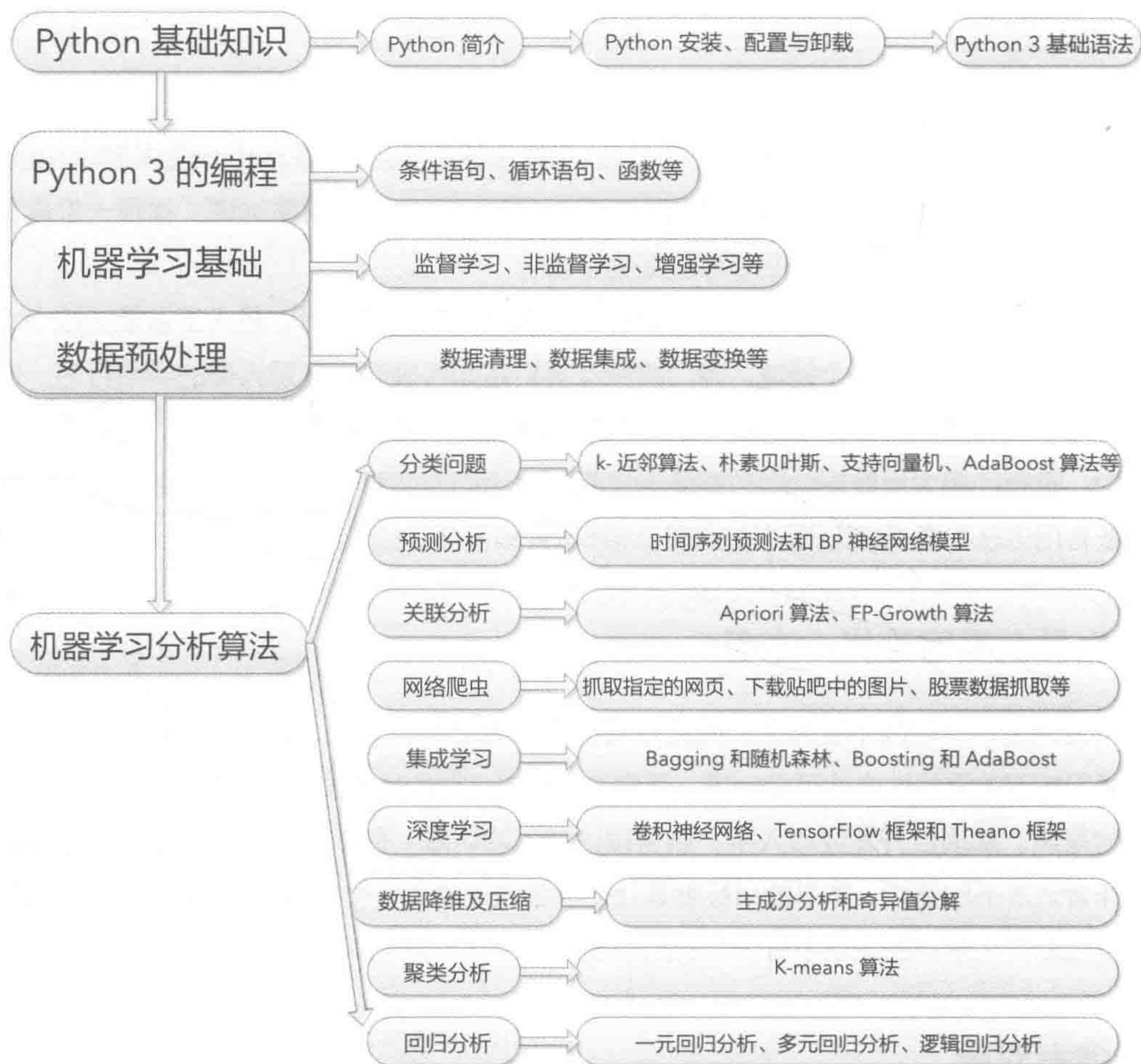
◆ 为什么要写这样一本书

古人云:“不闻不若闻之,闻之不若见之,见之不若知之,知之不若行之。”实践对于学习知识的重要性由此可见一斑,而理论知识与实践经验的脱节,恰恰是很多 Python 图书的写照。从项目开发经验入手,结合理论知识的讲解,便成为本书的立足点,也转化成了作者对本书的要求。本书的目标就是让初学者快速成为一个 Python 程序员,并拥有项目开发技能,在未来的职场中有一个较高的起点。

◆ 读者对象

- 没有任何 Python 基础的初学者
- 有一定的 Python 基础,想精通 Python 数据分析与机器学习的人员
- 有一定的 Python 基础,没有项目经验的人员
- 正在进行毕业设计的学生
- 大专院校及培训学校的教师和学生

◆ Python 最佳学习途径



◆ 本书特色

■ 零基础也能入门

无论您是否从事计算机相关行业，是否接触过 Python，是否使用 Python 开发过项目，都能从本书开启学习之旅。

■ 专业的项目指导

本书结合实际工作中的范例逐一讲解 Python 的各种知识和技术，使您在实战中掌握知识，轻松拥有项目经验。

◆ 配套资源

■ 源代码及视频下载

扫描下方二维码或在浏览器中输入下载链接“<http://v.51pcbook.cn/download/29566.html>”，即可下载本书配套素材与视频资源。



① **提示：**如果下载链接失效，请加入“办公之家”QQ群（218192911），联系管理员获取最新下载链接。

① **注意：**如加入QQ群时，系统提示此群已满，请根据验证信息加入新群。

■ 扫描二维码观看同步视频

使用微信、QQ或浏览器中的“扫一扫”功能，扫描每节中对应的二维码，即可观看相应的同步教学视频。

■ 手机版打包视频

读者可以扫描下方二维码下载龙马高新教育手机APP，将其直接安装到手机中，随时随地问同学、问专家，尽享海量资源。同时，我们也会不定期地向读者推送学习中常见的难点、使用技巧、行业应用等精彩内容，让学习更加简单高效。



◆ 作者团队

本书由龙马高新教育策划，由史卫亚任主编，于俊伟任副主编。其中第0~9章、11、14、16章由河南工业大学史卫亚老师编著，第10、12、13、15章由河南工业大学于俊伟

老师编著。在编写过程中，编者竭尽所能地为读者呈现最好、最全的实用功能，但仍难免有疏漏和不妥之处，敬请广大读者指正。读者若在学习过程中产生疑问，或有任何建议，可以通过以下方式联系我们。

投稿邮箱：pup7@pup.cn

读者信箱：2751801073@qq.com

读者交流 QQ 群：218192911（办公之家）

注意：如加入 QQ 群时，系统提示此群已满，请根据验证信息加入新群。

目 录


第 0 章 本书的技术体系	1
0.1 Python 的发展趋势.....	2
0.2 人工智能时代学习 Python 的重要性.....	2
0.3 本书的技术体系.....	3
0.4 学习本书需要注意的事项.....	6
第 1 章 Python 基础知识	7
1.1 Python 简介.....	8
1.1.1 了解 Python 的起源与发展历史.....	8
1.1.2 Python 的特色.....	8
1.1.3 学习 Python 的原因.....	9
1.2 Python 的当前版本.....	9
1.3 Python 的优缺点.....	10
1.4 Python 与其他语言的区别.....	10
1.5 Python 的应用领域.....	11
第 2 章 Python 的安装、配置与卸载	13
2.1 Python 的安装.....	14
2.1.1 Python 的下载.....	14
2.1.2 Python 的安装.....	15

2.2	Python 的配置.....	17
2.2.1	Python 环境变量的设置.....	17
2.2.2	Python 的启动.....	18
2.3	Python 的卸载.....	19

第 3 章 Python 3 基础语法.....21

3.1	第一个 Python 程序.....	22
3.2	Python 的输入和输出.....	26
3.2.1	Python 的输出语句.....	26
3.2.2	Python 的输入语句.....	26
3.3	Python 的基本数据类型.....	27
3.3.1	数字.....	27
3.3.2	字符串.....	28
3.3.3	列表.....	29
3.3.4	元组.....	30
3.3.5	集合.....	31
3.3.6	字典.....	32
3.4	Python 库的导入.....	32
3.5	Python 的集成开发环境.....	34
 3.6	自测练习.....	36


第 4 章 Python 3 的编程.....37

4.1	条件语句.....	38
4.2	循环语句.....	39
4.2.1	while 循环.....	40
4.2.2	for 循环.....	40
4.3	函数.....	43
4.4	模块.....	45
 4.5	自测练习.....	46


第 5 章	机器学习基础	47
5.1	机器学习概述	48
5.2	监督学习简介	50
5.3	非监督学习简介	50
5.4	增强学习简介	51
5.5	深度学习简介	53
5.6	机器学习常用术语	55
第 6 章	Python 机器学习及分析工具	57
6.1	矩阵操作函数库 (NumPy)	58
6.1.1	NumPy 的安装	58
6.1.2	NumPy 的基本使用	59
6.2	科学计算的核心包 (SciPy)	64
6.2.1	科学计算的核心包的安装	64
6.2.2	科学计算的核心包的基本使用	67
6.3	Python 的绘图库 (Matplotlib)	71
6.3.1	Matplotlib 简介及安装	71
6.3.2	Matplotlib 的基本使用	72
6.4	数据分析包 (Pandas)	79
6.4.1	Pandas 简介和安装	79
6.4.2	Pandas 的基本使用方法	80
6.5	机器学习函数库 (Scikit-learn)	81
6.6	统计建模工具包 (StatsModels)	88
6.7	深度学习框架 (TensorFlow)	90
第 7 章	数据预处理	93
7.1	数据预处理概述	94
7.2	数据清理	95

7.2.1	异常数据处理	95
7.2.2	缺失值处理	96
7.2.3	噪声数据处理	97
7.3	数据集成.....	98
7.4	数据变换.....	99
7.5	数据归约.....	101
7.6	Python 的主要数据预处理函数.....	102
7.6.1	Python 的数据结构	102
7.6.2	数据缺失处理函数	104

第 8 章 分类问题..... 115

8.1	分类概述.....	116
8.2	常用方法.....	116
8.2.1	k- 近邻算法	116
8.2.2	朴素贝叶斯	118
8.2.3	支持向量机	123
8.2.4	AdaBoost 算法	125
8.2.5	决策树	127
8.2.6	Multi-layer Perceptron 多层感知机.....	135
8.3	项目实战.....	137
8.3.1	实例 1: 使用 k- 近邻算法实现约会网站的配对效果	137
8.3.2	实例 2: 使用朴素贝叶斯过滤垃圾邮件	141
8.3.3	实例 3: SVM 实现手写识别系统.....	144
8.3.4	实例 4: 基于单层决策树构建分类算法	147
8.3.5	实例 5: 使用决策树对 iris 数据集分类	151
8.3.6	实例 6: 使用决策树对身高体重数据进行分类	153
8.3.7	实例 7: 使用 k- 近邻算法对鸢尾花数据进行交叉验证	156
8.3.8	使用多层感知器分析, 根据葡萄酒的各项化学特征来 判断葡萄酒的优劣	160
 8.4	自测练习.....	163

第 9 章 预测分析.....165


- 9.1 预测概述..... 166
- 9.2 常用方法..... 166
 - 9.2.1 时间序列分析预测法 166
 - 9.2.2 BP 神经网络模型 168
- 9.3 项目实战..... 170
 - 9.3.1 实例 1: 根据一年的历史数据预测后十年的数据趋势 170
 - 9.3.2 实例 2: 使用神经网络预测公路运量 178
-  9.4 自测练习.....185

第 10 章 关联分析.....187


- 10.1 关联分析概述..... 188
- 10.2 基本方法..... 188
 - 10.2.1 Apriori 算法..... 189
 - 10.2.2 FP-Growth 算法 189
- 10.3 项目实战 (解决目前流行的实际问题) 192
 - 10.3.1 用 Apriori 进行关联分析的实例..... 192
 - 10.3.2 使用 FP-Growth 算法提取频繁项集 195
-  10.4 自测练习.....199

第 11 章 网络爬虫.....201

- 11.1 网络爬虫概述..... 202
 - 11.1.1 网络爬虫原理..... 202
 - 11.1.2 爬虫分类..... 203
- 11.2 网页抓取策略和方法 204
 - 11.2.1 网页抓取策略..... 204
 - 11.2.2 网页抓取的方法..... 204
- 11.3 项目实战 205

11.3.1	用 Python 抓取指定的网页	205
11.3.2	用 Python 抓取包含关键词的网页	207
11.3.3	下载贴吧中的图片	208
11.3.4	股票数据抓取	210
 11.4	自测练习	213


第 12 章 集成学习.....215

12.1	集成学习概述	216
12.2	常用方法	216
12.2.1	Bagging 和随机森林	216
12.2.2	Boosting 和 AdaBoost	217
12.3	项目实战	219
12.3.1	使用随机森林方法预测乘员的存活概率	219
12.3.2	使用 AdaBoost 方法进行二元分类	222
 12.4	自测练习	225


第 13 章 深度学习.....227

13.1	深度学习概述	228
13.2	常用方法	228
13.2.1	监督学习的深度学习网络结构	229
13.2.2	非监督学习的深度学习网络结构	230
13.3	项目实战	233
13.3.1	使用 TensorFlow 框架进行 MNIST 数据集生成	233
13.3.2	使用 Theano 框架进行 MNIST 数字识别	237
 13.4	自测练习	241

第 14 章 数据降维及压缩.....243

- 14.1 数据降维及压缩概述 244
 - 14.1.1 数据降维 244
 - 14.1.2 图像压缩 245
- 14.2 基本方法..... 245
 - 14.2.1 主成分分析 245
 - 14.2.2 奇异值分解 248
- 14.3 项目实战..... 251
 - 14.3.1 主成分分析 PCA 实例..... 251
 - 14.3.2 使用奇异值分解进行图像压缩 257
-  14.4 自测练习.....260

第 15 章 聚类分析.....261

- 15.1 聚类分析概述..... 262
- 15.2 K-means 算法 264
 - 15.2.1 K-means 算法与步骤 264
 - 15.2.2 K-means 算法涉及的问题 264
 - 15.2.3 实际聚类问题的处理流程 265
- 15.3 项目实战..... 266
 - 15.3.1 K-means 算法实现二维数据聚类 266
 - 15.3.2 使用 Scikit-learn 中的方法进行聚类分析 269
-  15.4 自测练习.....276

第 16 章 回归分析问题279

- 16.1 回归分析概述..... 280
- 16.2 基本方法..... 280
 - 16.2.1 一元回归分析 280
 - 16.2.2 多元线性回归 281

16.2.3	回归的计算方法	282
16.2.4	逻辑回归分析	284
16.3	项目实战.....	286
16.3.1	身高与体重的回归分析	286
16.3.2	房价预测	293
16.3.3	产品销量与广告的多元回归分析	296
16.3.4	鸢尾花数据的逻辑回归分析	298
 16.4	自测练习.....	300

0

第 0 章 本书的技术体系

Python 语言是一种语法简洁而清晰，面向对象的高级程序设计语言，并且 Python 语言使用广泛，代码范例也很多，便于读者快速学习和掌握，十分容易上手。机器学习已应用于人们生活的方方面面，远超出大多数人的想象，如人脸识别、预测天气、垃圾邮件过滤和电商网站购物产品推荐等。随着各种数据以指数级增长，我们不仅需要使用更好的工具对这些数据进行分析，还要通过这些数据分析掌握其内涵，进而进行学习以提高人类应对未知世界的能力。

本章将介绍以下内容：

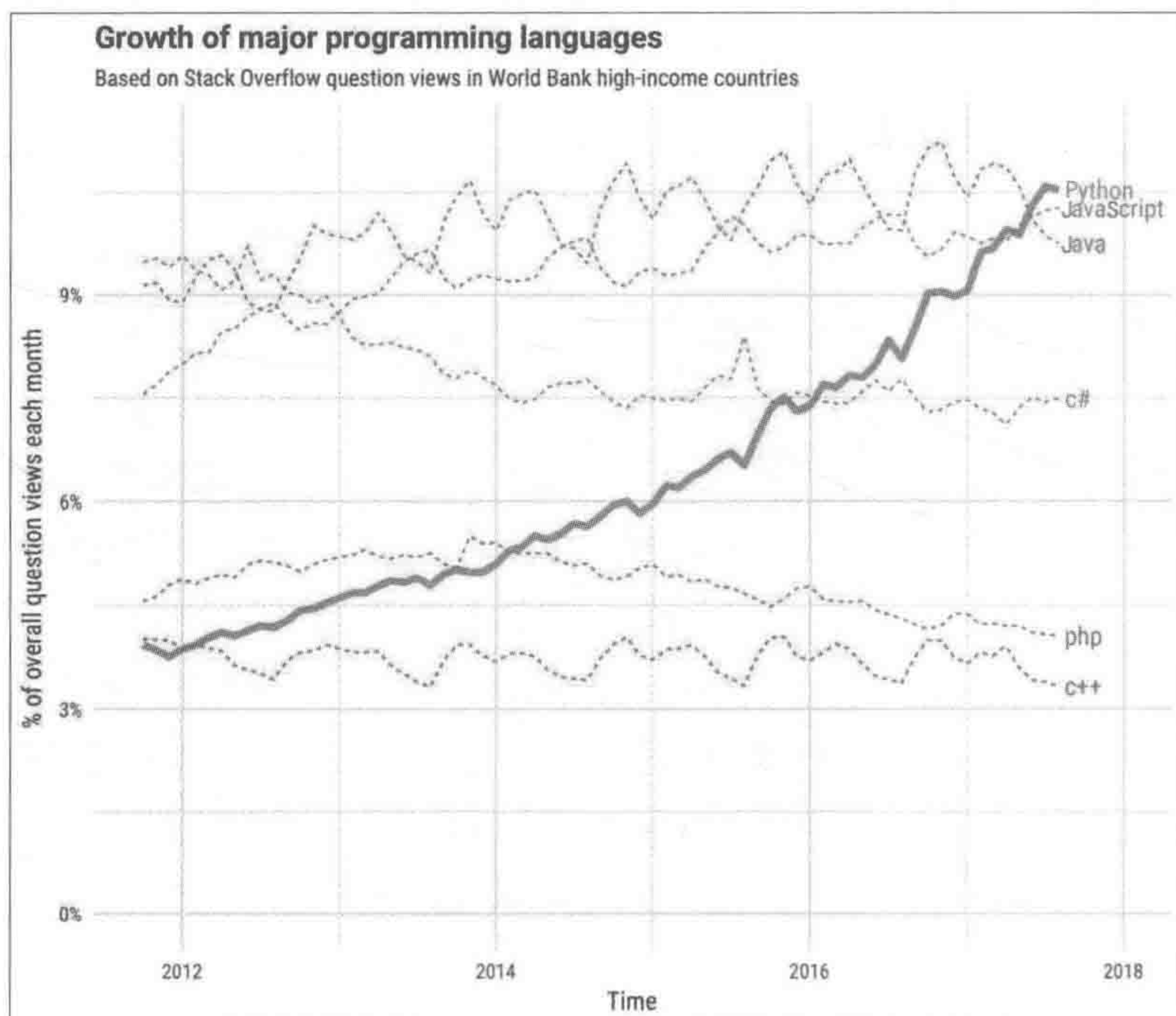
- Python 的发展趋势
- 人工智能时代学习 Python 的重要性
- 本书的技术体系
- 学习本书需要注意的事项

0.1 Python 的发展趋势



也许几年前提到 Python 语言，很多人知之甚少，然而最近一两年不管是程序员，还是普通的计算机爱好者，都逐渐接触到这门编程语言，Python 的关注量也占据榜首。

知名 IT 技术问答社区 Stack Overflow 最近公布了程序语言排行榜，排名结果如下图所示。



从 Stack Overflow 编程语言流行趋势中可以看到，Python 的热度在过去几年中一直在迅速增长，已经成为目前热度增长最快的编程语言。2017 年 6 月，Python 第一次成为高收入国家在 Stack Overflow 中访问量最大的编程语言。

因为 Python 语言的语法非常简单易懂，所以很多提及编程就恐慌的人少了一些担心，Python 语言不仅在学术上非常受欢迎，而且很多非计算机专业的人也在学习 Python。他们利用自己写的简单的小程序，让生活变得精彩起来。

0.2 人工智能时代学习 Python 的重要性



现在的社会是一个高速发展的社会，科技发达，信息流通，人们之间的