

采用Python 3.6版本，兼容Python 3.X等众多版本

# Python

## 数据分析从入门到精通

、Numpy、Matplotlib、  
数据分析专家。

张啸宇 李静 编著

# Python

## 数据分析从入门到精通

张啸宇 李静 编著

电子工业出版社

Publishing House of Electronics Industry

北京·BEIJING

## 内 容 简 介

对于希望使用 Python 来完成数据分析工作的人来说,学习 IPython、Numpy、pandas、Matplotlib 这个组合是目前看来不错的方向。本书就是这样一本循序渐进的书。

本书共 3 篇 14 章。第 1 篇是 Python 数据分析语法入门,将数据分析用到的一些语言的语法基础讲解清楚,为接下来的数据分析做铺垫。第 2 篇是 Python 数据分析工具入门,介绍了 Python 数据分析“四剑客”——IPython、Numpy、pandas、Matplotlib。第 3 篇是 Python 数据分析案例实战,包括两个案例,分别是数据挖掘和玩转大数据,为读者能真正使用 Python 进行数据分析奠定基础。

本书内容精练、重点突出、实例丰富,是广大数据分析工作者必备的参考书,同时也非常适合大、中专院校师生学习阅读,还可作为高等院校统计分析及相关专业的教材。

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

Python 数据分析从入门到精通 / 张啸宇,李静编著. —北京:电子工业出版社,2018.3

ISBN 978-7-121-33613-3

I. ①P… II. ①张… ②李… III. ①软件工具—程序设计 IV. ①TP311.561

中国版本图书馆 CIP 数据核字(2018)第 020314 号

策划编辑:张月萍

责任编辑:牛 勇

特约编辑:赵树刚

印 刷:三河市双峰印刷装订有限公司

装 订:三河市双峰印刷装订有限公司

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:787×980 1/16 印张:20.75 字数:415 千字

版 次:2018 年 3 月第 1 版

印 次:2018 年 6 月第 2 次印刷

印 数:1000 册 定价:69.00 元

凡所购买电子工业出版社图书有缺损问题,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, 88258888。

质量投诉请发邮件至 [zltz@phei.com.cn](mailto:zltz@phei.com.cn), 盗版侵权举报请发邮件到 [dbqq@phei.com.cn](mailto:dbqq@phei.com.cn)。

本书咨询联系方式:010-51260888-819, [faq@phei.com.cn](mailto:faq@phei.com.cn)。

# 前言

---

由于 Python 具有简单、易学、免费开源、可移植性、可扩展性等特点，所以它的受欢迎程度扶摇直上。再加上 Python 拥有非常丰富的库，这也使得它在数据分析领域有着越来越广泛的应用。如果你已经决定学习 Python 数据分析，但是之前没有编程经验，那么本书将会是你的正确选择。

本书的第 1 篇主要介绍学习数据分析必备的一些 Python 语法基础，包括 Python 的安装、数据类型、数据结构、模块、类、异常处理、使用 pip 安装 Python 需要的一些工具等；第 2 篇主要介绍 Python 在数据处理和科学计算方面的工具和方法，包括 IPython 交互式壳的使用、Jupyter Notebook 的使用和 Numpy 的使用，还介绍了 Python 的核心数据分析处理库 pandas，以及 Python 著名的 2D 绘图库 Matplotlib；第 3 篇通过数据挖掘和玩转大数据两个案例总结和应用前面所学的知识。

这三篇的层进正好是 Python 数据分析入门者的阶梯，读者通过学习这三部分内容，即可迈入数据分析的门槛。

## 本书的特点

Python 是当前非常流行的面向对象编程语言，本书将其在数据分析处理方面的特色发挥到极致。本书的主要特点如下：

- Python 被大量应用在数据挖掘和机器学习领域，其中使用极其广泛的是 IPython、Numpy、pandas、Matplotlib 等库。本书详细地介绍了这些库的组成与使用，为科学计算相关人员提供了有用的参考资料。
- 本书采取循序渐进的写作风格，对于工具的安装、使用步骤、方法技巧逐步展

开，加以图解和应用场景，即使完全不懂 Python 和数据分析的人员，也可以流畅地读完本书。

- 无论哪种语言，编程的方法、模式、数据结构、算法都是相通的。本书将科学计算、数据结构与各种工具和方法完美结合，让非 Python 读者也能融会贯通，让学习统计的人能找到更适合的统计方法和数据分析处理方法。
- 本书最后的两个实战案例适合数据分析入门者，案例的步骤详细、分析到位，能为读者入手真实项目打下良好的基础。

## 本书的内容安排

本书共 3 篇 14 章，主要章节规划如下：

第 1 章介绍了 Python 的发展历程、特性，帮助读者搭建最基础的数据分析环境，下载开发语言，选择开发工具，然后在此基础上开发自己的第一个 Python 程序。读者在学完本章内容后应该对 Python 有一个基础的认识，知道为什么选择它来进行数据分析。

第 2 章介绍了 Python 的基础语法，包括它的代码组织形式、如何缩进、如何注释等，以及输入/输出该如何处理，在中文环境下如何更好地使用 Python 是本章的重点，最后还通过一个实例复习了 Python 的这些语法。读者在学完本章内容后可以轻松地编写一些简单的 Python 程序。

第 3 章介绍了 Python 的数据类型与流程控制语句。如果读者已有编程基础，那么阅读本章内容不会有任何压力。如果没有编程基础，那么学习一门语言的流程控制最关键的就是这些知识。读者在阅读完本章后就能轻松阅读更大的 Python 程序。

第 4 章介绍了可复用的函数与模块。这些内容较为复杂，但却是进行数据分析的关键。每个数据处理过程我们都会用到函数或模块，而我们后期用到的数据分析库也可以说是一个大函数。所以学习完本章内容，读者应该能够看明白一个完整的 Python 库。

第 5 章介绍了数据结构与算法，这是数据分析的基础，也是人工智能的基础。利用算法我们可以找到解决方案，也可以找到最优路径，还可以更高效地完成数据分析任务。读者如果没有看懂本章内容，一定要反复阅读，直到学会为止。

第 6 章介绍了面向对象的 Python。面向对象已经成为每门语言都具备的特性，类、对象、继承这些概念都是面向对象的基础。如果读者没有编程经验，则阅读本章可能会

有一定的难度，但是了解了对象的概念，就能学会如何编写更高效的代码、如何让代码和代码之间联动起来。

第 7 章介绍了异常处理与程序调试。机器毕竟不是人，如果出现错误，则可能会导致死机，或者数据出错。为了防止这些错误的发生，或者防止程序的使用方能得到反馈，我们必须学会 Python 的异常处理功能。并且当程序发生错误时，我们还要通过程序调试找到错误所在。

第 8 章介绍了 pip 软件包管理。既然在做数据分析时我们要用到很多数据分析库，那么如何下载、安装或管理这些库就成了数据分析的第一步。pip 就是这样一个工具，它能下载、安装、更新、显示、搜索我们需要的数据分析库。

第 9 章介绍了 IPython 科学计算库，它是使用 Python 进行数据分析、处理、呈现的重要选择之一。本章主要介绍了 Python 科学计算库的安装方法、IPython 壳的一些特性和基本功能、Jupyter Notebook 的安装和使用方法。IPython 壳的使用是本章的重点，也是数据分析处理的基础工具，希望读者能够消化本章内容，为真正做好数据项目打下基础。

第 10 章介绍了 Numpy 科学计算库，主要介绍了它的数组对象及数组对象的一些基本属性和生成数组的基本方法，还包括数组的索引和分片等基本操作，这部分内容是 Numpy 数据处理的核心。本章介绍的代数运算函数、线性代数、统计函数等内容会让读者觉得有些困惑，但这已经进入了数据分析的关键时刻，所以仍建议读者对本章的内容融会贯通。

第 11 章介绍了 pandas 数据分析处理库，主要包括它的序列、数据框的基本操作，还包括 pandas 里处理默认值、读取常见格式的文本数据，以及数据的组合和分组操作。最后介绍了 pandas 的时间序列和一个处理实际数据集的案例，读者掌握了这些内容，就可以更好地处理数据。

第 12 章介绍了 Matplotlib 的 Pyplot 和 Artist 模块，以及 pandas 的绘图功能。对于读者来说，Pyplot 模块是需要掌握的，Artist 模块是需要了解的，pandas 的绘图功能在实际数据分析中要能熟练应用。

第 13 章是数据挖掘的案例。首先介绍了著名的贝叶斯理论，然后实现了贝叶斯分类器，最后实现了协同过滤算法，这些都是数据挖掘、分析领域的基础算法。建议读者尝试自己编写代码，熟练掌握贝叶斯分类器和协同过滤算法的使用。

第 14 章是玩转大数据的案例。鉴于本书主要针对数据分析入门者，所以本章也逐

步实现了数据的分析过程，从了解数据到分析数据，最后到代码实现，相信读者学完本章内容后，就能真正动手分析大数据了。

本书由浅入深、从理论到实践，尤其适合初学者逐步学习和完善自己的知识结构。

请访问 [www.broadview.com.cn/33613](http://www.broadview.com.cn/33613) 下载本书配套代码。

## 适合阅读本书的读者

- 希望从事数据分析相关工作的人员。
- 数据分析工作人员。
- 大数据从业人员。
- Python 爱好者。
- 人工智能从业人员。
- 统计行业的人员。
- 大、中专院校统计相关专业的学生。

编者

# 目 录

## 第 1 篇 Python 数据分析语法入门

第 1 章 初识 Python .....	1
1.1 Python 是什么 .....	2
1.2 Python 有什么优点 .....	3
1.2.1 Python 是自由开源的软件 .....	3
1.2.2 Python 是跨平台的 .....	3
1.2.3 Python 功能强大 .....	4
1.2.4 Python 是可扩展的 .....	4
1.2.5 Python 易学易用 .....	5
1.3 其他程序设计语言中的 Python .....	5
1.3.1 Jython.....	5
1.3.2 Python for .NET.....	6
1.3.3 IronPython.....	6
1.4 快速搭建 Python 开发环境 .....	7
1.4.1 Python 的下载和安装 .....	7
1.4.2 用 Visual Studio 编译 Python 源代码.....	9
1.4.3 Python 开发工具: Vim.....	10
1.4.4 Python 开发工具: Emacs.....	15
1.4.5 Python 开发工具: PythonWin .....	18
1.4.6 其他的 Python 开发工具 .....	20



1.5	第一个 Python 程序 .....	22
1.5.1	从“Hello, Python!”开始 .....	22
1.5.2	Python 的交互式命令行 .....	24
1.6	本章小结.....	25
<b>第 2 章</b>	<b>Python 起步必备 .....</b>	<b>27</b>
2.1	Python 代码的组织形式和注释方式 .....	27
2.1.1	用缩进来分层.....	28
2.1.2	代码的两种注释方式.....	29
2.1.3	Python 语句的断行 .....	30
2.2	Python 的基本输入/输出函数 .....	31
2.2.1	接收输入的 input()函数.....	31
2.2.2	输出内容的 print()函数.....	32
2.3	Python 对中文的支持 .....	33
2.3.1	Python 3 之前的版本如何使用中文.....	33
2.3.2	更全面的中文支持.....	36
2.4	简单实用的 Python 计算器 .....	37
2.4.1	直接进行算术运算.....	37
2.4.2	math 模块提供丰富的数学函数.....	38
2.4.3	Python 对大整数的支持 .....	39
2.5	本章小结.....	40
<b>第 3 章</b>	<b>Python 的数据类型与流程控制语句 .....</b>	<b>41</b>
3.1	Python 数据类型：数字 .....	42
3.1.1	整型和浮点型.....	42
3.1.2	运算符.....	43
3.2	Python 数据类型：字符串 .....	45
3.2.1	Python 中的字符串 .....	45
3.2.2	字符串中的转义字符.....	46

3.2.3	操作字符串 .....	46
3.2.4	字符串的索引和分片 .....	49
3.2.5	格式化字符串 .....	50
3.2.6	字符串、数字类型的转换 .....	50
3.2.7	原始字符串 .....	51
3.3	Python 数据类型：列表和元组 .....	52
3.3.1	创建和操作列表 .....	52
3.3.2	创建和操作元组 .....	53
3.4	Python 数据类型：字典 .....	54
3.5	Python 数据类型：文件 .....	55
3.6	Python 数据类型：布尔值 .....	56
3.7	Python 的流程控制语句 .....	56
3.7.1	分支结构：if 语句 .....	57
3.7.2	循环结构：for 语句 .....	59
3.7.3	循环结构：while 语句 .....	62
3.8	本章小结 .....	63
<b>第 4 章</b>	<b>可复用的函数与模块 .....</b>	<b>64</b>
4.1	Python 自定义函数 .....	65
4.1.1	函数的定义 .....	65
4.1.2	函数调用 .....	66
4.2	参数让函数更有价值 .....	67
4.2.1	有默认值的参数 .....	67
4.2.2	参数的传递方式 .....	69
4.2.3	如何传递任意数量的参数 .....	70
4.2.4	用参数返回计算结果 .....	70
4.3	变量的作用域 .....	71
4.4	最简单的函数：使用 lambda 表达式定义函数 .....	72

4.5	可重用结构: Python 模块 .....	73
4.5.1	Python 模块的基本用法 .....	73
4.5.2	Python 在哪里查找模块 .....	75
4.5.3	是否需要编译模块 .....	77
4.5.4	模块也可独立运行 .....	78
4.5.5	如何查看模块提供的函数名 .....	79
4.6	用包来管理多个模块 .....	80
4.6.1	包的组成 .....	80
4.6.2	包的内部引用 .....	81
4.7	本章小结 .....	81
<b>第 5 章</b>	<b>数据结构与算法 .....</b>	<b>82</b>
5.1	表、栈和队列 .....	82
5.1.1	表 .....	83
5.1.2	栈 .....	84
5.1.3	队列 .....	86
5.2	树和图 .....	88
5.2.1	树 .....	88
5.2.2	二叉树 .....	89
5.2.3	图 .....	93
5.3	查找与排序 .....	95
5.3.1	查找 .....	96
5.3.2	排序 .....	97
5.4	本章小结 .....	100
<b>第 6 章</b>	<b>面向对象的 Python .....</b>	<b>101</b>
6.1	面向对象编程概述 .....	101
6.1.1	Python 中的面向对象思想 .....	102
6.1.2	类和对象 .....	102

6.2	在 Python 中定义和使用类 .....	103
6.2.1	类的定义 .....	104
6.2.2	类的使用 .....	105
6.3	类的属性和方法 .....	106
6.3.1	类的属性 .....	107
6.3.2	类的方法 .....	108
6.4	类的继承 .....	111
6.4.1	使用继承 .....	111
6.4.2	Python 的多重继承 .....	112
6.5	在类中重载方法和运算符 .....	114
6.5.1	方法重载 .....	114
6.5.2	运算符重载 .....	115
6.6	在模块中定义类 .....	117
6.7	本章小结 .....	119
<b>第 7 章</b>	<b>异常处理与程序调试 .....</b>	<b>120</b>
7.1	异常的处理 .....	120
7.1.1	使用 try 语句捕获异常 .....	121
7.1.2	常见异常的处理 .....	123
7.1.3	多重异常的捕获 .....	124
7.2	用代码引发异常 .....	125
7.2.1	使用 raise 语句引发异常 .....	126
7.2.2	assert——简化的 raise 语句 .....	127
7.2.3	自定义异常类 .....	128
7.3	使用 pdb 模块调试 Python 脚本 .....	128
7.3.1	调试语句块 .....	129
7.3.2	调试表达式 .....	129
7.3.3	调试函数 .....	130
7.3.4	设置断点 .....	131
7.3.5	pdb 调试命令 .....	131

7.4	在 PythonWin 中调试脚本.....	134
7.5	本章小结.....	136
<b>第 8 章</b>	<b>pip 软件包管理.....</b>	<b>137</b>
8.1	安装 pip .....	137
8.2	更新 pip .....	138
8.3	pip 常用操作 .....	138
8.3.1	安装软件包.....	138
8.3.2	卸载软件包.....	139
8.3.3	更新软件包.....	139
8.3.4	显示本地所有已经安装的软件包.....	139
8.3.5	显示软件包的细节.....	139
8.3.6	搜索软件包.....	140
8.3.7	通过 wheel 文件安装软件包 .....	141
8.4	本章小结.....	141

## 第 2 篇 Python 数据分析工具入门

<b>第 9 章</b>	<b>IPython 科学计算库 .....</b>	<b>142</b>
9.1	IPython 简介 .....	143
9.2	安装 IPython 及其他相关库 .....	144
9.2.1	使用 Anaconda 安装 .....	144
9.2.2	使用 pip 安装.....	145
9.3	IPython 壳基础.....	146
9.3.1	自动补全.....	147
9.3.2	检查.....	149
9.3.3	%run 命令 .....	150
9.3.4	快捷键.....	150
9.3.5	异常和错误定位.....	151

9.3.6	魔法方法.....	151
9.3.7	和操作系统交互.....	152
9.3.8	代码分析: %prun 和%run .....	153
9.3.9	目录标签系统.....	155
9.3.10	嵌入 IPython.....	155
9.4	融合 Matplotlib 库和 Pylab 模型.....	156
9.5	输入和输出变量.....	157
9.6	交互式调试器.....	158
9.7	计时功能.....	159
9.8	重新载入模块.....	160
9.9	配置 IPython.....	161
9.10	Jupyter.....	162
9.10.1	基于 Qt 的控制台 .....	162
9.10.2	Jupyter Notebook .....	165
9.11	IPython 和 Jupyter Notebook 的关系.....	170
9.12	本章小结.....	173
<b>第 10 章</b>	<b>Numpy 科学计算库 .....</b>	<b>174</b>
10.1	Numpy 基础.....	174
10.1.1	数组对象介绍.....	175
10.1.2	生成数组.....	176
10.1.3	数组对象数据类型.....	180
10.1.4	打印数组.....	182
10.2	数组的基本操作.....	184
10.3	基本的分片和索引操作.....	186
10.4	高级索引.....	189
10.4.1	整数索引.....	189
10.4.2	布尔索引.....	190
10.4.3	布尔索引的简单应用.....	192

10.5	改变数组的形状.....	193
10.6	组装、分割数组.....	195
10.7	数组的基本函数.....	196
10.8	复制和指代.....	198
10.9	线性代数.....	199
10.10	使用数组来处理数据.....	201
10.11	Numpy 的 where()函数和统计函数.....	203
10.11.1	where()函数.....	203
10.11.2	统计函数.....	205
10.12	输入与输出.....	206
10.12.1	二进制文件.....	206
10.12.2	文本文件.....	207
10.13	生成随机数.....	208
10.14	数组的排序和查找.....	210
10.14.1	排序.....	210
10.14.2	查找.....	212
10.15	扩充转换.....	213
10.16	本章小结.....	215
<b>第 11 章</b>	<b>pandas 数据分析处理库.....</b>	<b>216</b>
11.1	pandas 数据结构介绍.....	217
11.1.1	序列.....	217
11.1.2	数据框.....	221
11.2	索引对象.....	226
11.3	核心的基本函数.....	227
11.4	索引和旋转.....	229
11.5	算术运算与对齐.....	232
11.6	处理默认值.....	233
11.7	多级索引.....	237

11.8	读/写数据.....	239
11.9	组合数据.....	243
11.10	数据分组操作.....	247
11.11	时间序列.....	249
11.11.1	时间序列介绍.....	250
11.11.2	使用时间序列作图.....	253
11.12	本章小结.....	259
<b>第 12 章</b>	<b>Matplotlib 数据可视化.....</b>	<b>260</b>
12.1	Pyplot 模块介绍.....	261
12.1.1	plot()函数.....	261
12.1.2	绘制子图.....	264
12.1.3	添加注释.....	266
12.1.4	其他的坐标轴类型.....	268
12.2	应用 Pyplot 模块.....	269
12.3	Artist 模块.....	275
12.3.1	Artist 模块概述.....	275
12.3.2	Artist 的属性.....	277
12.4	使用 pandas 绘图.....	283
12.5	本章小结.....	287

### 第 3 篇 Python 数据分析案例实战

<b>第 13 章</b>	<b>案例 1: 数据挖掘.....</b>	<b>288</b>
13.1	贝叶斯理论介绍.....	288
13.2	贝叶斯分类器的实现.....	290
13.3	协同过滤推荐系统.....	295
13.3.1	相似度计算.....	296
13.3.2	协同过滤推荐系统的实现.....	300
13.4	本章小结.....	304



第 14 章 案例 2: 玩转大数据 .....	305
14.1 案例概述 .....	306
14.1.1 了解大数据的处理方式 .....	306
14.1.2 处理日志文件 .....	307
14.1.3 案例目标 .....	308
14.2 日志文件的分割 .....	309
14.3 编写 Map() 函数处理小文件 .....	311
14.4 编写 Reduce() 函数 .....	313
14.5 本章小结 .....	315