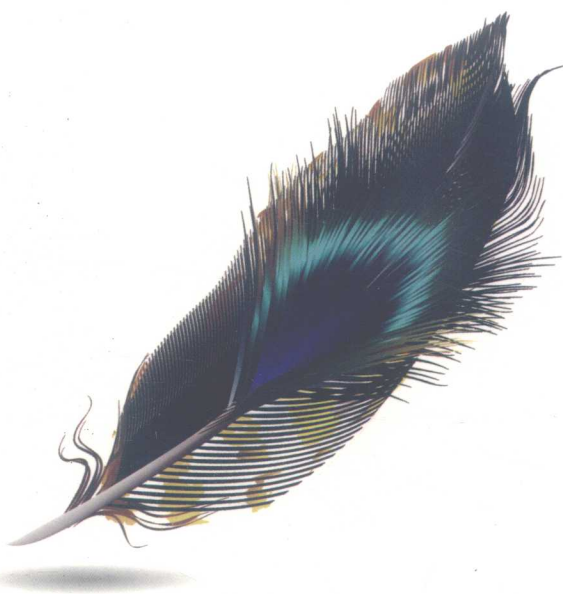




浙江省普通高校“十三五”新形态教材

大数据与人工智能技术丛书



# Python 数据分析与实践

◎ 柳毅 主编 毛峰 李艺 副主编

400分钟  
视频讲解

教学课件

教学大纲

程序源码

电子教案

习题答案

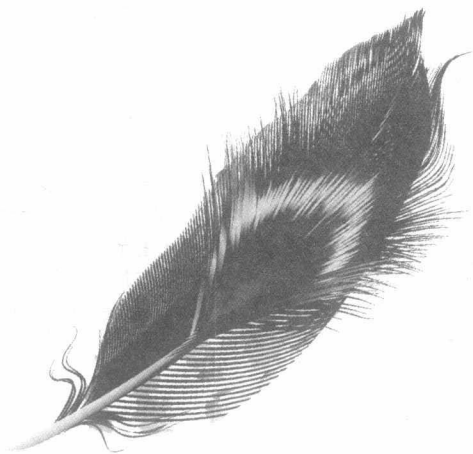
清华大学出版社





浙江省普通高校“十三五”新形态教材

大数据与人工智能技术丛书



# Python 数据分析与实践

◎ 柳毅 主编 毛峰 李艺 副主编

清华大学出版社

北京

## 内 容 简 介

Python 是信息管理与信息系统、电子商务等信息管理类本科学生进行数据分析所需要掌握的基础性语言和分析工具,是未来学生掌握大数据分析技术的学习基础。本书共分 12 章,着重讲述 Python 语言 and 数据分析工具包的应用。第 1 章主要介绍 Python 的发展历史、特点、集成开发环境、内置模块、帮助的使用等内容;第 2 章主要介绍 Python 语言的基础知识;第 3 章主要介绍 Python 中的常用数据结构,包括序列、字典、集合等,以及函数的定义和调用等;第 4 章主要介绍 Python 中类、对象和方法的相关内容;第 5 章主要介绍 Python 进行数据分析常用的 NumPy、Pandas、Matplotlib、SciPy 和 Scikit-learn 等基础库内容;第 6 章主要介绍网络数据获取的 HTML 和 XML 两种网页组织形式,以及 urllib 和 BeautifulSoup4 两个模块内容;第 7 章主要介绍文件的操作;第 8 章主要介绍数据可视化,以及使用 Python 绘制图表的知识;第 9 章主要介绍利用 Python 进行数据库应用开发;第 10、11 章主要介绍 Python 机器学习的基本概念以及有监督、无监督学习算法的原理;第 12 章主要介绍 Python 在地理空间分析上的应用。本书中的代码均在 Python 3.5 中测试通过。

本书一方面侧重对 Python 数据分析基础知识的讲解,另一方面注重 Python 数据处理方法的应用。本书适合作为计算机科学与技术专业学生学习数据分析的入门教材,也适合作为 Python 爱好者的参考书。

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 数据分析与实践/柳毅主编. —北京:清华大学出版社,2019  
(大数据与人工智能技术丛书)  
ISBN 978-7-302-51579-1

I. ①P… II. ①柳… 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 257419 号

策划编辑:魏江江  
责任编辑:王冰飞  
封面设计:刘 键  
责任校对:时翠兰  
责任印制:刘海龙

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](mailto:c-service@tup.tsinghua.edu.cn)

质量反馈:010-62772015, [zhiliang@tup.tsinghua.edu.cn](mailto:zhiliang@tup.tsinghua.edu.cn)

课件下载: <http://www.tup.com.cn>, 010-62795954

印 装 者:清华大学印刷厂

经 销:全国新华书店

开 本:185mm×260mm

印 张:21.5

字 数:370 千字

版 次:2019 年 6 月第 1 版

印 次:2019 年 6 月第 1 次印刷

印 数:1~1500

定 价:59.00 元

产品编号:080337-01

# 序

人类社会已经进入数字经济时代,大数据、云计算、机器学习、人工智能等技术纷至沓来,数据的管理和应用已经渗透到每一个行业的业务领域,成为当今乃至将来企业运作的基础资产。只有掌握数据并善于运用数据的人,才会在未来社会日益激烈的竞争环境中保持领先地位。

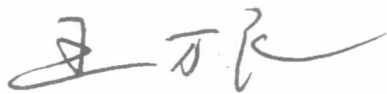
Python 语言很好地融合了大数据分析、机器学习以及人工智能技术,是目前大数据和机器学习领域热门的语言之一。本书为学习者深入浅出地介绍 Python 数据分析的原理、建模过程、统计应用方法,具有极强的实践性。

本书基于 Python 3.5 工具环境,通过实践案例讲解 Python 控制、处理、分析数据的算法和工具,让学生了解如何利用 Python 编程和数据处理库(包括 NumPy、SciPy、Matplotlib、Pandas 及 Scikit-learn 等)高效地解决各种数据分析问题,发挥 Python 在数据分析、可视化、机器学习、地理空间信息分析方面的优势,引导读者成为数据分析的高手。

本书内容严谨,逻辑清晰,可供计算机科学与技术专业以及信息管理与信息系统、电子商务等信息管理类的本科与研究生学习使用,为大数据时代的企业管埋、市场营销、金融等行业从事大量数据分析的从业人员提供科学的学习资源。

浙江工业大学计算机科学与技术学院院长

国家级教学名师、教授、博士生导师



2019年3月

# 前言

Python 是大数据时代非常受欢迎的数据分析编程语言,近年来随着机器学习、云计算、人工智能等技术的发展,Python 的流行趋势扶摇直上,已经成为数据分析和数据科学事实上的标准语言和标准平台之一。

本书针对数据分析人员和 Python 编程学习者进行内容编排和章节讲述,Python 数据分析整个学习路线图计划分成 16 周,120 天左右。本书主要内容包括以下五大部分。

(1) Python 工作环境及基础语法知识:认识 Python 程序运行方式,使用 Python 3.5 开发集成环境与工具;学习 Python 程序基本结构,理解 Python 的面向对象定义和类、对象的操作方法,以及 Python 异常处理机制。本部分为基础内容,建议学习时间为 4 周。

(2) Python 数据分析相关知识:Python 生态系统为分析师和数据科学家提供了各种程序库,例如 NumPy、SciPy、Pandas 和 Matplotlib,使其在数据分析领域也有广泛的应用。Python 数据分析的学习主要是对相关库的使用,例如数据整理需要用到 NumPy 库,数据描述与分析则主要用到 Pandas 库。由于有前面的学习基础,本部分学习时间建议为 3 周。

(3) Python 数据可视化:Python 数据可视化的过程就是学习 Matplotlib 库的过程,Matplotlib 库包含丰富的数据可视化资源,地图、3D 等都有涉及,基于前面两部分的学习经验,这部分内容在两周内基本可以完成。

(4) Python 机器学习:Scikit-learn 是本书所使用的核心程序库,依托于上述几种工具包,封装了大量的经典以及最新的机器学习模型。通过介绍有监督和无监督机器学习原理,学习有监督学习的线性回归、Logistic 回归、朴素贝叶斯、SVM、KNN 和决策树等几个常用算法,以及无监督学习的 K-Means 聚类算法。在前面三部分学习的基础上,本部分内容建议学习时间在 4 周左右。

(5) Python 地理空间数据分析应用:向读者介绍地理空间分析的基本概念和常

用的地理空间数据,然后介绍 Python 中与地理数据处理和地理分析相关的工具,最后以 Python 处理和分析矢量数据与栅格数据的方法,对浙江省实施的“五水共治”行动中劣五类水体在地理空间模型上分布的卫星图像数据进行可视化数据分析的综合应用,本部分内容建议学习时间为 3 周。

本书编写人员具有丰富的 Python 数据分析实践经验和多年的信息管理教学能力,第 1~3 章由沈阳工业大学李艺老师编写;第 4~7、10 章由杭州电子科技大学柳毅老师编写;第 8、9、11、12 章由杭州电子科技大学毛峰老师编写;王健、陆佳涣等硕士研究生参与了本书相关章节内容和程序代码的完善工作;浙江工业大学计算机科学与技术学院院长、国家教学名师王万良教授对本书进行了认真的审阅,并提出许多宝贵的建设性意见,使本书内容日臻完善,在此对他们所付出的辛勤劳动表示诚挚的感谢。

本书结合大数据管理与应用的最新发展,针对计算机科学与技术、信息管理与信息系统、电子商务等经管类本科教学特点进行撰写。本书提供教学课件、教学大纲、电子教案、习题答案和程序源码,读者可以扫描封底的课件二维码下载。本书还提供 400 分钟的视频讲解,扫描书中的二维码,可以在线观看。

由于编者水平所限,书中难免有疏漏之处,敬请读者批评指正。

编者

2019 年 3 月

# 目 录








源码下载

<b>第 1 章 Python 简介</b> .....	1
1.1 Python 语言的发展史 .....	1
1.1.1 Python 语言的特点 .....	4
1.1.2 Python 2 与 Python 3 的区别 .....	6
1.2 Python 的环境搭建  .....	7
1.3 开始使用 Python IDLE  .....	10
1.3.1 交互方式 .....	10
1.3.2 Python 的集成开发环境 .....	11
1.4 Eclipse+PyDev 的安装 .....	14
1.5 代码风格 .....	20
1.6 使用帮助 .....	26
本章小结 .....	28
习题 .....	28
<b>第 2 章 Python 语言基础知识</b> .....	29
2.1 标识符与变量 .....	29
2.1.1 标识符 .....	29
2.1.2 变量 .....	30
2.2 数据类型及运算 .....	33
2.2.1 数据类型  .....	34
2.2.2 运算符和表达式  .....	35
2.3 分支结构控制语句  .....	39
2.3.1 if 语句 .....	39
2.3.2 if-else 语句 .....	40
2.3.3 if-elif-else 语句 .....	41

2.4 循环语句 .....	42
2.4.1 循环结构控制语句  .....	42
2.4.2 循环嵌套控制语句 .....	43
2.4.3 break 语句和 continue 语句 .....	43
2.4.4 range()函数 .....	45
2.5 常见的 Python 函数 .....	46
本章小结 .....	52
习题 .....	52
<b>第3章 数据结构与函数设计</b> .....	<b>53</b>
3.1 序列 .....	53
3.1.1 列表  .....	54
3.1.2 元组 .....	56
3.1.3 字符串 .....	57
3.1.4 列表与元组之间的转换 .....	59
3.2 字典  .....	59
3.2.1 创建字典 .....	59
3.2.2 字典的方法 .....	60
3.2.3 列表、元组与字典之间的转换 .....	60
3.3 集合 .....	61
3.3.1 集合的创建 .....	61
3.3.2 集合的运算 .....	63
3.3.3 集合的方法 .....	65
3.4 函数的定义 .....	67
3.4.1 函数的调用  .....	69
3.4.2 形参与实参 .....	69
3.4.3 函数的返回  .....	70
3.4.4 位置参数  .....	70
3.4.5 默认参数与关键字参数  .....	71
3.4.6 可变长度参数 .....	72
本章小结 .....	73











习题 .....	74
<b>第 4 章 类与对象</b> .....	<b>75</b>
4.1 面向对象 .....	75
4.1.1 面向对象编程  .....	76
4.1.2 类的抽象与封装  .....	77
4.2 认识 Python 中的类、对象和方法 .....	78
4.2.1 类的定义与创建  .....	78
4.2.2 构造函数 .....	81
4.3 类的属性 .....	82
4.3.1 类属性和实例属性 .....	82
4.3.2 公有属性和私有属性 .....	83
4.4 类的方法 .....	85
4.4.1 类方法的调用 .....	85
4.4.2 类方法的分类  .....	85
4.4.3 析构函数 .....	87
4.5 类的继承 .....	88
4.5.1 父类与子类 .....	88
4.5.2 继承的语法  .....	88
4.5.3 多重继承 .....	90
4.5.4 运算符的重载 .....	92
4.6 类的组合 .....	93
4.7 类的异常处理 .....	97
4.7.1 异常 .....	97
4.7.2 Python 中的异常类 .....	98
4.7.3 捕获与处理异常 .....	99
4.7.4 自定义异常类 .....	103
4.7.5 with 语句 .....	104
4.7.6 断言 .....	105
本章小结 .....	107
习题 .....	107

案例 .....	108
<b>第 5 章 Python 数据分析基础库</b> .....	110
5.1 NumPy .....	111
5.1.1 ndarray 的数据类型  .....	113
5.1.2 数组和标量之间的运算  .....	114
5.1.3 索引和切片  .....	114
5.1.4 数组转置和轴对换  .....	117
5.1.5 利用数组进行数据处理  .....	118
5.1.6 数学和统计方法 .....	120
5.2 Pandas .....	121
5.2.1 Pandas 数据结构 .....	121
5.2.2 Pandas 文件操作 .....	123
5.2.3 数据处理 .....	124
5.2.4 层次化索引 .....	125
5.2.5 分级顺序 .....	128
5.2.6 使用 DataFrame 的列 .....	129
5.3 Matplotlib .....	130
5.3.1 figure 和 subplot .....	131
5.3.2 调整 subplot 周围的间距 .....	134
5.3.3 颜色、标记和线型 .....	135
5.3.4 刻度标签和图例 .....	135
5.3.5 添加图例 .....	136
5.3.6 将图表保存到文件 .....	137
5.4 SciPy .....	138
5.5 Scikit-learn .....	139
本章小结 .....	141
习题 .....	141
<b>第 6 章 网络数据的获取</b> .....	142
6.1 网页数据的组织形式 .....	143
6.1.1 HTML .....	143

6.1.2	HTML 元素	144
6.1.3	HTML 属性	146
6.2	XML	147
6.2.1	XML 的结构和语法	148
6.2.2	XML 元素和属性	150
6.3	利用 urllib 处理 HTTP 	153
6.4	利用 BeautifulSoup4 解析 HTML 文档	158
6.4.1	BeautifulSoup4 中的对象 	160
6.4.2	遍历文档树 	163
6.4.3	搜索文档树 	168
	本章小结	177
	习题	177
<b>第 7 章</b>	<b>文件操作</b>	<b>178</b>
7.1	文件的打开和关闭 	178
7.1.1	打开文件	178
7.1.2	关闭文件	180
7.2	读写文件 	180
7.2.1	从文件读取数据	180
7.2.2	向文件写入数据	182
7.3	文件对话框	182
7.3.1	基于 win32ui 构建文件对话框	182
7.3.2	基于 tkinter 构建文件对话框	183
7.4	应用实例：文本文件的操作	184
	本章小结	188
	习题	189
<b>第 8 章</b>	<b>Python 数据可视化</b>	<b>190</b>
8.1	数据可视化概念框架	190
8.1.1	数据可视化简介	190
8.1.2	数据可视化常用图表	192
8.1.3	Python 数据可视化环境准备	195

8.2 绘制图表 .....	197
8.2.1 Matplotlib API 入门  .....	197
8.2.2 创建图表  .....	198
8.2.3 图表定制  .....	204
8.2.4 保存图表 .....	208
8.3 更多高级图表及定制 .....	208
8.3.1 样式 .....	208
8.3.2 subplot 子区  .....	210
8.3.3 图表颜色和填充 .....	212
8.3.4 动画 .....	213
本章小结 .....	215
习题 .....	215
<b>第9章 数据库应用开发</b> .....	<b>216</b>
9.1 Python 与数据库 .....	216
9.1.1 数据库简介 .....	216
9.1.2 Python 数据库工作环境 .....	220
9.2 本地数据库 SQLite  .....	223
9.2.1 SQLite 简介 .....	223
9.2.2 Python 内置的 sqlite3 模块 .....	223
9.3 关系型数据库  .....	225
9.3.1 关系型数据库基本操作与 SQL .....	225
9.3.2 操作 MySQL .....	226
9.4 非关系型数据库 .....	232
9.4.1 NoSQL 介绍 .....	232
9.4.2 MongoDB  .....	234
9.4.3 PyMongo: MongoDB 和 Python  .....	236
习题 .....	241
<b>第10章 机器学习——有监督学习</b> .....	<b>242</b>
10.1 机器学习简介 .....	242
10.2 Python 机器学习库 Scikit-learn .....	243

10.3 有监督学习 .....	245
10.3.1 线性回归  .....	246
10.3.2 Logistic 回归分类器 .....	248
10.3.3 朴素贝叶斯分类器 .....	252
10.3.4 支持向量机 .....	257
10.3.5 KNN 算法  .....	259
10.3.6 决策树  .....	264
本章小结 .....	272
习题 .....	272
<b>第 11 章 机器学习——无监督学习</b> .....	<b>273</b>
11.1 无监督学习 .....	273
11.2 聚类 .....	274
11.2.1 相异度 .....	274
11.2.2 K-Means 算法  .....	277
11.2.3 DBSCAN 算法  .....	282
11.3 关联规则 .....	286
11.3.1 关联分析 .....	286
11.3.2 Apriori 算法  .....	288
11.3.3 FP-growth 算法 .....	294
本章小结 .....	303
习题 .....	303
<b>第 12 章 Python 地理空间分析</b> .....	<b>304</b>
12.1 地理空间分析简介 .....	304
12.1.1 地理空间分析的基本概念 .....	304
12.1.2 地理空间分析与 Python .....	305
12.2 地理空间数据 .....	306
12.2.1 数据格式概览 .....	306
12.2.2 数据特征 .....	307
12.2.3 矢量数据 .....	307
12.2.4 栅格数据 .....	309

12.3 Python 地理空间分析工具	309
12.3.1 GeoJSON	309
12.3.2 GDAL 和 OGR	311
12.3.3 PyShp	311
12.3.4 PIL	312
12.3.5 GeoPandas	313
12.4 Python 分析矢量数据 	313
12.4.1 访问矢量数据	313
12.4.2 Shapefile 文件操作	314
12.4.3 空间查询	315
12.4.4 叠加分析	316
12.5 Python 与遥感 	317
12.5.1 访问影像文件	317
12.5.2 影像裁剪	318
12.5.3 重采样	321
12.5.4 影像分类	321
12.6 “五水共治”资源地理空间分析综合应用	323
本章小结	327
习题	327

# 第 1 章



## Python简介

---

本章学习目标：

- 了解 Python 语言的发展历史及特点
- 深刻了解 Python 2 和 Python 3 的区别
- 熟练掌握 Python 中 IDLE 的编程特点
- 熟练掌握 Python 的开发环境

本章首先向读者介绍 Python 发展的历史及其特点；然后分析 Python 2 与 Python 3 不同的编程特点；接着以 Eclipse+PyDev 为例对 Python 的环境搭建进行介绍，并对 Python 自带的 IDLE 界面进行讲解；最后对 Python 的各种开发环境进行讲解。

### 1.1 Python 语言的发展史

Python 的作者是荷兰人 Guido von Rossum，尽管拥有阿姆斯特丹大学数学和计算机双硕士学位，Guido 总趋向于做计算机相关的工作，并热衷于做任何与编程相关的活儿。Guido 接触并使用过 Pascal、C、Fortran 等语言，这些语言的基本设计原则是让计算机能更快地运行。在 20 世纪 80 年代，虽然 IBM 和苹果公司已经掀起了个

人计算机浪潮,但这些个人计算机的配置很低。所有的编译器的核心是做优化,以便让程序能够运行。为了提高效率,语言也迫使程序员像计算机一样思考,以便能写出更适合计算机的程序。在那个时代,程序员恨不得拥有使用计算机每一点空间的能力,有人甚至认为 C 语言的指针是在浪费内存。至于动态类型、内存自动管理、面向对象等,那就不用想了,否则会让计算机陷入瘫痪。

这种编程方式让 Guido 感到苦恼。Guido 知道如何用 C 语言写出一个功能,但整个编写过程需要耗费大量的时间。他的另一个选择是 Shell。Bourne Shell 作为 UNIX 系统的解释器已经长期存在,UNIX 的管理员们经常用 Shell 去写一些简单的脚本,以进行一些系统维护的工作,比如定期备份、文件系统管理等。许多用 C 语言编写上百行的程序,在 Shell 下只用几行就可以完成。然而,Shell 的本质是调用命令,它并不是一个真正的语言。比如说,Shell 没有数值型的数据类型,即使加法运算也很复杂。总之,Shell 不能全面地调用计算机的功能。

Guido 希望有一种语言能够像 C 语言那样可以全面调用计算机的功能接口,又能够像 Shell 那样可以轻松地编程,ABC 语言让 Guido 看到希望。ABC 是由荷兰的数学和计算机研究所开发的。Guido 在该研究所工作,并参与到 ABC 语言的开发中。与当时的大部分语言不同,ABC 语言的目标是“让用户感觉更好”。ABC 语言希望让程序变得容易阅读、使用、记忆和学习,并以此来激发人们学习编程的兴趣。尽管已经具备了良好的可读性和易用性,但 ABC 语言最终没有流行起来。在当时 ABC 语言的设计也存在一些致命的问题。

(1) ABC 语言不是模块化语言。如果想在 ABC 语言中增加功能,比如对图形化的支持,就必须改动很多地方。

(2) ABC 语言不能直接操作文件系统。尽管用户可以通过诸如文本流的方式导入数据,但 ABC 语言无法直接读写文件,输入输出的困难对于计算机语言来说是致命的。

(3) ABC 语言用自然语言的方式来表达程序的意义,然而对于程序员来说,他们更习惯用 function 或者 define 来定义一个函数,用等号来分配变量。尽管 ABC 语言很特别,但学习难度也很大。

(4) ABC 语言编译器很大,必须被保存在磁带上,这样 ABC 语言很难快速传播。

1989 年,为了打发圣诞节假期,Guido 开始写 Python 语言的编译器。“Python”这个名字来自 Guido 所挚爱的电视剧 *Monty Python's Flying Circus*。他希望



Python 语言能符合他的理想——创造一种 C 和 Shell 之间的功能全面、易学易用、可拓展的语言。1991 年,第一个 Python 编译器诞生。它是用 C 语言实现的,并能够调用 C 语言的库文件。从一诞生,Python 已经具有了类、函数、异常处理,包含表和词典在内的核心数据类型,以及以模块为基础的拓展系统。

Python 语法很多来自 C 语言,但又受到 ABC 语言的很大影响。来自 ABC 语言的一些规定直到今天还有争议,比如强制缩进,但这些语法规则让 Python 容易阅读。另外,Python 聪明地选择服从一些惯例,特别是 C 语言的惯例,比如回归等号赋值。

Python 从一开始就特别注重可拓展性,Python 可以在多个层次上拓展。用户可以直接引入 .py 文件,也可以引用 C 语言的库。Python 程序员可以快速地使用 Python 写 .py 文件作为拓展模块,但当性能是考虑的重要因素时,Python 程序员可以深入底层写 C 程序,编译为 .so 文件引入到 Python 中使用。Python 就好像是使用钢结构建房一样,先规定好大的框架,而程序员可以在此框架下自由地拓展或更改。

最初的 Python 完全由 Guido 本人开发。Python 受到 Guido 同事的欢迎,他们迅速地反馈使用意见,并参与到 Python 的改进中。Guido 和一些同事构成 Python 的核心团队,他们将自己的大部分业余时间用于 hack Python。随后,Python 被拓展到研究所之外。Python 将许多机器层面上的细节隐藏,交给编译器处理,并凸显出逻辑层面的编程思考。Python 程序员可以花更多时间用于思考程序的逻辑,而不是具体的实现细节,这一特征吸引了广大的程序员,Python 开始流行。

Guido 维护了一个邮件列表,Python 用户能通过邮件进行交流。Python 用户来自许多领域,不同的背景对 Python 也有不同的需求。Python 相当开放,又容易拓展,所以当用户不满足现有功能时很容易对 Python 进行拓展或改造。随后,这些用户将改动发给 Guido,并由 Guido 决定是否将新的特征加入到 Python 或者标准库中。如果代码能被纳入 Python 自身或者标准库,这将是极大的荣誉。由于 Guido 有着至高无上的决定权,所以他被称为“终身的仁慈独裁者”。

Python 被称为“Battery Included”,是说它的标准库功能强大。这是整个社区的贡献,Python 的开发者来自不同领域,他们将不同领域的优点带给 Python,比如 Python 标准库中的正则表达式参考 Perl,而 lambda 匿名函数以及 map()、filter()、reduce() 等函数参考了 Lisp。在 Python 的开发过程中,社区起到了重要的作用。