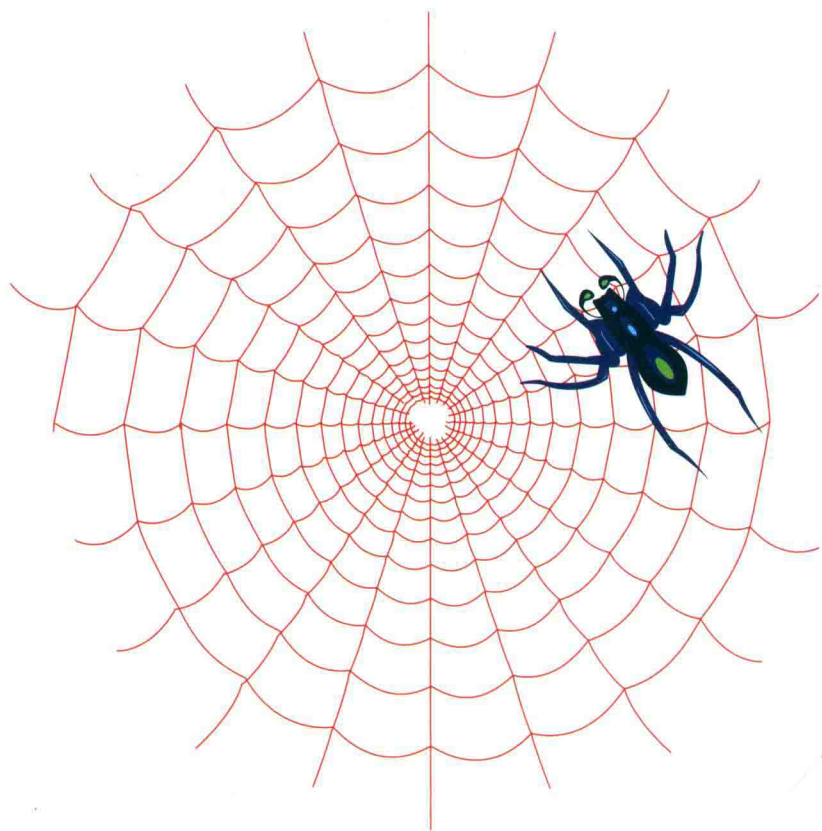


基于Python 3.7，快速掌握网络爬虫编程

- 针对想直接切入Python 3网络爬虫编程的读者
- 围绕爬虫技术、爬虫实战，提高你的数据获取能力
- 学习过程中贯穿大小示例，配合教学视频，便于读者掌握爬虫编程技巧

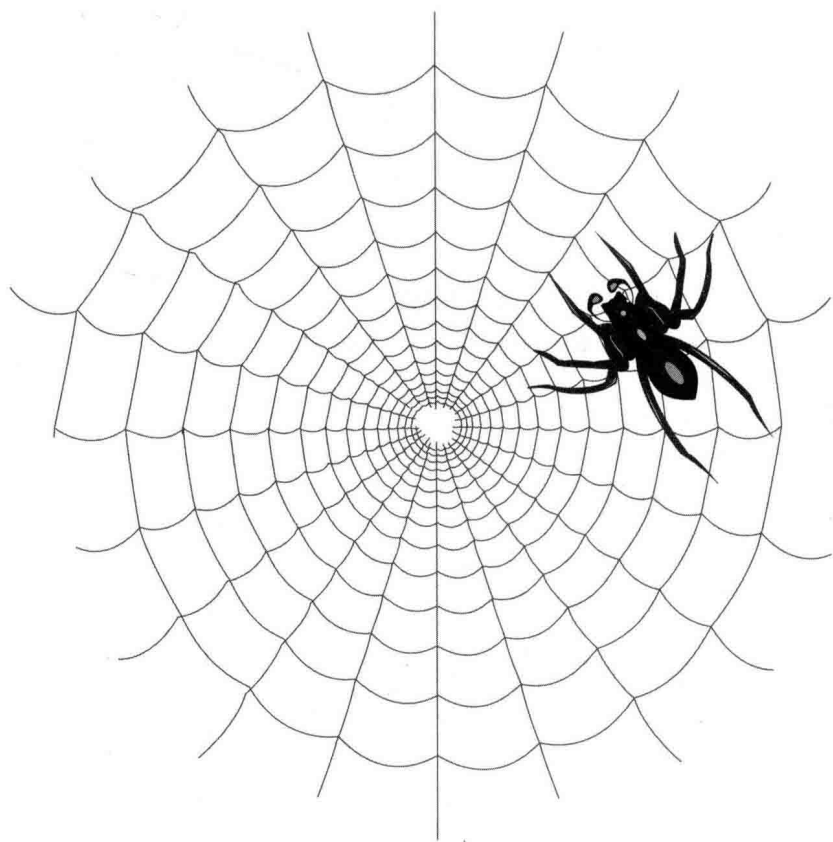


Python 3.7网络爬虫 快速入门

王启明 编著



清华大学出版社



Python 3.7网络爬虫 快速入门

王启明 编著

清华大学出版社
北京

内 容 简 介

Python 3.7 正在成为目前流行的编程语言，而网络爬虫又是 Python 网络应用中的重要技术，二者的碰撞产生了巨大的火花。本书在这个背景下编写而成，详细介绍 Python 3.7 网络爬虫技术。

本书分为 11 章，分别介绍 Python 3.7 爬虫开发相关的基础知识、lxml 模块、BeautifulSoup 模块、正则表达式、文件处理、多线程爬虫、图形识别、Scrapy 框架、PyQuery 模块等。基本上每一章都配有众多小范例程序与一个大实战案例。作者还为每一章分别录制教学视频供读者自学参考。

本书内容详尽、示例丰富，是有志于学习 Python 网络爬虫技术初学者必备的参考书，同时也可作为 Python 爱好者拓宽知识领域、提升编程技术的参考书。

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 3.7 网络爬虫快速入门/王启明编著. — 北京：清华大学出版社，2019
ISBN 978-7-302-53647-5

I. ①P… II. ①王… 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 (2019) 第 186648 号

责任编辑：夏毓彦

封面设计：王 翔

责任校对：闫秀华

责任印制：丛怀宇

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：北京鑫海金澳胶印有限公司

经 销：全国新华书店

开 本：190mm×260mm 印 张：13.25 字 数：339 千字

版 次：2019 年 10 月第 1 版 印 次：2019 年 10 月第 1 次印刷

定 价：49.00 元

产品编号：082462-01

前言

Python 是简练的语言

使用像 Python 这样的动态类型语言编写的代码往往比用其他主流语言编写的代码更加简短。这意味着，在编程的过程中会有更少的录入工作，而且会更容易记住算法并真正领会算法的原理。

Python 是易读的语言

Python 不时被人们指为“可执行的伪代码”。虽然很明显这是夸大之词，但是它表明大多数有经验的程序员可以读懂 Python 代码并领会代码所要表达的意图。

Python 是易安装的语言

要搭建 Python 的环境非常容易，不管是 Windows、Linux 还是 Mac 系统，只要配置好 Python 的环境，只需要 `easy_install XX` 或者 `pip install XX` 就可以安装所需要的第三方扩展包。

Python 是易扩展的语言

Python 附属了很多标准库，涉及数据函数、XML 解析以及网页下载、RSS 解析、SQLite 等，可以解决现实中遇到的各种问题。

为什么用 Python 实现网络爬虫

基于上述优点，加上抓取网页文档的接口更简洁；相比其他动态脚本语言，如 Perl、Shell，Python 的 `urllib2` 包提供了较为完整的访问网页文档的 API，以及抓取后的处理方法，比如筛选 HTML 标签、提取文本等。Python 的相关扩展可以用极短的代码完成大部分文档的处理。

本书涉及的技术或框架

Python 基本语法

Python 函数

lxml 模块

XPath 语法

BeautifulSoup

正则表达式

XML

CSV

MySQL

PyQuery

线程 (Thread)

进程 (Process)

图形识别验证码

Scrapy

本书涉及的范例和案例

爬取豆瓣网的内容

爬取电影天堂网的内容

爬取猫眼电影网的内容

爬取腾讯招聘网的内容

使用 BeautifulSoup 爬取电影天堂的内容

使用正则表达式爬取糗事百科的内容

爬取鼠绘漫画的图片

使用多线程爬取豆瓣的内容

使用 Tesseract 识别验证码的图片

使用 Scrapy 框架爬取豆瓣网的内容

使用 PyQuery 爬取微博热搜的内容

本书特点

(1) 语言通俗易懂。对于没有基础的读者，最忌讳的就是讲一些艰深晦涩的理论，让人难于理解。本书则尽量使用通俗易懂的语言来介绍 Python，让大家更容易理解各种知识点，从而将相应的知识变成自己的能力。

(2) 结合范例程序来讲解知识点。为了讲明各个知识点，基本上每个知识点都通过相关的范例程序来说明。通过范例程序及实际的执行效果，让大家学以致用，在理解领会的基础上进一步掌握相关知识、相应模块的方法。

(3) 插图配合教学视频。为了保证本书的范例程序均能够成功执行，每个范例程序不仅有相应的程序代码，还有程序执行后的效果图。大家可以通过效果图来对比程序执行的结果，确保学习质量。同时每章还配有一个教学视频供读者自学参考。

(4) 案例丰富。为了向读者说明 Python 爬虫程序的效果，书中选择的被爬取的网站都是国内热门的网站，比如豆瓣电影、猫眼电影、电影天堂、微博热搜等。这些网站大家喜闻乐见。通过这些案例，大家可以轻松地掌握相关模块的使用方法，举一反三，将相应技术应用于其他同类的网站中。

代码与教学视频下载

本书示例源代码与教学视频下载地址请扫描右边二维码获得。

如果下载有问题，请联系 booksaga@163.com，邮件主题为“Python 3.7 网络爬虫快速入门”。



本书读者

- 有志于学习 Python 爬虫编程的初学者
- 对 Python 网络爬虫技术有兴趣的开发人员
- 各类综合信息网站的站长或技术人员
- 高校和培训学校相关专业的师生

编者
2019年7月

目 录

第 1 章 简识 Python	1
1.1 了解 Python.....	1
1.1.1 Python 的概念.....	1
1.1.2 有趣的 Python 程序.....	2
1.2 集成开发环境.....	4
1.2.1 安装 Python 3.7.....	4
1.2.2 从 IDLE 启动 Python.....	6
1.3 编写自己的第一个 Python 程序：一个简单的问候.....	8
1.4 小结.....	11
第 2 章 Python 语法速览.....	12
2.1 数据类型与变量.....	12
2.1.1 数据类型.....	12
2.1.2 变量.....	14
2.2 运算符.....	15
2.2.1 算术运算符.....	16
2.2.2 比较运算符.....	17
2.2.3 赋值运算符.....	17
2.2.4 逻辑运算符.....	18
2.2.5 位运算符.....	19
2.2.6 成员运算符.....	20
2.2.7 身份运算符.....	21
2.2.8 运算符的优先级.....	21
2.3 使用复合类型.....	21
2.3.1 列表.....	22
2.3.2 元组.....	26
2.3.3 字典.....	26

2.3.4 集合	27
2.4 流程控制结构	29
2.4.1 选择结构	29
2.4.2 重复结构（循环结构）	30
2.5 小结	33
第3章 函数	34
3.1 认识函数	34
3.1.1 什么是函数	34
3.1.2 创建函数	35
3.2 使用函数	35
3.2.1 参数	36
3.2.2 返回值	38
3.2.3 函数的递归	39
3.3 实践一下	40
3.3.1 实践一：编写一个函数	40
3.3.2 实践二：遍历与计数	41
3.4 小结	42
第4章 lxml 模块和 XPath 语法	43
4.1 lxml 模块	43
4.1.1 什么是模块	43
4.1.2 关于 lxml 模块	44
4.1.3 lxml 模块的安装	44
4.1.4 lxml 库的用法	46
4.2 XPath 语法	46
4.2.1 基本语法	46
4.2.2 基本操作	47
4.2.3 lxml 库的用法	49
4.2.4 XPath 范例程序测试	50
4.3 爬虫 lxml 解析实战	53
4.3.1 爬取豆瓣网站	53
4.3.2 爬取电影天堂	55
4.3.3 爬取猫眼电影	58
4.3.4 爬取腾讯招聘网	61
4.3.5 关于 HTML	63
4.4 小结	63

第 5 章 BeautifulSoup 库	64	
5.1 简识 BeautifulSoup 4	64	
5.1.1 安装与配置	64	
5.1.2 基本用法	66	
5.2 BeautifulSoup 对象	67	
5.2.1 创建 BeautifulSoup 对象	67	
5.2.2 4 类对象	70	
5.2.3 遍历文档树	74	
5.2.4 搜索文档树	78	
5.3 方法和 CSS 选择器	81	
5.3.1 find 类方法	81	
5.3.2 CSS 选择器	82	
5.4 爬取示范：使用 BeautifulSoup 爬取电影天堂	85	
5.4.1 基本思路	85	
5.4.2 实际爬取	85	
5.5 小结	87	
第 6 章 正则表达式	88	
6.1 了解正则表达式	88	
6.1.1 基本概念	88	
6.1.2 re 模块	89	
6.1.3 compile()方法	89	
6.1.4 match()方法	90	
6.1.5 group()和 groups()方法	90	
6.1.6 search()方法	90	
6.1.7 findall()方法	92	
6.1.8 finditer()方法	93	
6.1.9 split()方法	94	
6.1.10 sub()方法	94	
6.2 抓取	95	
6.2.1 抓取标签间的内容	95	
6.2.2 抓取 tr <td> 标签间的内容</td> <td>98</td>	标签间的内容	98
6.2.3 抓取标签中的参数	99	
6.2.4 字符串处理及替换	101	
6.3 爬取实战	102	
6.3.1 获取数据	103	
6.3.2 筛选数据	104	

6.3.3	保存数据	107
6.3.4	显示数据	107
6.4	总结	108
第 7 章	JSON 文件处理、CSV 文件处理和 MySQL 数据库操作	109
7.1	简识 JSON	109
7.1.1	什么是 JSON	109
7.1.2	字典和列表转 JSON	110
7.1.3	将 JSON 数据转储到文件中	111
7.1.4	将一个 JSON 字符串加载为 Python 对象	111
7.1.5	从文件中读取 JSON	112
7.2	CSV 文件处理	113
7.2.1	读取 CSV 文件	113
7.2.2	把数据写入 CSV 文件	114
7.2.3	练习	115
7.3	MySQL 数据库	117
7.3.1	MySQL 数据库的安装	117
7.3.2	安装 MySQL 模块	127
7.3.3	连接 MySQL	127
7.3.4	执行 SQL 语句	128
7.3.5	创建表	129
7.3.6	插入数据	130
7.3.7	查看数据	132
7.3.8	修改数据	133
7.3.9	删除数据	135
7.3.10	实践操作	136
7.4	小结	139
第 8 章	多线程爬虫	140
8.1	关于多线程	140
8.1.1	基本知识	140
8.1.2	多线程的适用范围	141
8.2	多线程的实现	142
8.2.1	使用 <code>_thread</code> 模块创建多线程	142
8.2.2	关于 <code>Threading</code> 模块	145
8.2.3	使用函数方式创建线程	146
8.2.4	传递可调用的类的实例来创建线程	148

8.2.5	派生子类并创建子类的实例	149
8.3	使用多进程	150
8.3.1	创建子进程	150
8.3.2	将进程定义为类	151
8.3.3	创建多个进程	152
8.4	爬取示范：多线程爬取豆瓣电影	153
8.4.1	使用多进程进行爬取	154
8.4.2	使用多线程进行爬取	156
8.5	小结	158
第 9 章	图形验证识别技术	159
9.1	图像识别开源库：Tesseract	159
9.1.1	安装 Tesseract	159
9.1.2	设置环境变量	164
9.1.3	验证安装	166
9.2	对网络验证码的识别	168
9.2.1	读取网络验证码并识别	168
9.2.2	对验证码进行转化	169
9.3	小结	170
第 10 章	Scrapy 框架	171
10.1	了解 Scrapy	171
10.1.1	Scrapy 框架概述	171
10.1.2	安装	173
10.2	开发 Scrapy 的过程	176
10.2.1	Scrapy 开发步骤	176
10.2.2	Scrapy 保存信息的格式	177
10.2.3	项目中各个文件的作用	178
10.3	爬虫范例	179
10.3.1	Scrapy 爬取美剧天堂	179
10.3.2	Scrapy 爬取豆瓣网	182
10.3.3	Scrapy 爬取豆瓣网 II	186
10.4	总结	189
第 11 章	PyQuery 模块	190
11.1	PyQuery 模块	190
11.1.1	什么是 PyQuery 模块	190

11.1.2	PyQuery 模块的安装.....	190
11.2	PyQuery 模块用法.....	191
11.2.1	使用字符串初始化 PyQuery 对象.....	191
11.2.2	使用文件初始化 PyQuery 对象.....	192
11.2.3	使用 URL 初始化 PyQuery 对象.....	193
11.3	CSS 筛选器的使用.....	194
11.3.1	基本 CSS 选择器.....	194
11.3.2	查找节点.....	195
11.3.3	遍历结果并输出.....	197
11.3.4	获取文本信息.....	198
11.4	爬虫 PyQuery 解析实战.....	200
11.4.1	爬取猫眼票房.....	200
11.4.2	爬取微博热搜.....	201
11.5	小结.....	202

第 1 章

◀ 简识Python ▶

Python 语言是一种面向对象的计算机程序设计语言，一经发行便受到众多计算机开发者的喜爱。经过多年的完善补充，Python 语言越发显现出它的优点，成为大数据技术人员热衷学习研究的热门语言。Python 简洁清晰的编程风格，易于和计算机其他应用领域巧妙结合，使得 Python 适应面非常广。

读者通过本章的学习，可以初步了解 Python 的语言魅力，以及如何将 Python 用来解决一些数学上的小问题。但是要想完全掌握好一门语言，则需要不断地使用，同时还要在使用过程中不断地积累编程技巧，最后达到融会贯通。

本章主要涉及的知识点有：

- 认识 Python: Python 是一种解释型、交互式、面向对象、对初学者友好的语言
- 编程: 通过编程实现和计算机的对话，用 Python 编程实现一些简单的功能
- 了解主流的开发环境: 借助完善的开发环境进行编程，事半功倍
- Hello, World! : 编写第一个小程序，并学会处理程序中可能出现的问题

1.1 了解 Python

本节首先介绍 Python 的基本概念，这些概念是学习和使用 Python 编程的前提。理解了 Python 的这些基础概念之后才能为学好 Python 编程打下坚实的基础。

1.1.1 Python 的概念

首先来看看 Python 的定义。Python 是一种解释型、交互式、面向对象的程序设计语言，也是对初学者友好的一种程序设计语言。

- Python 是解释型语言

Python 作为一种解释型语言，意味着在开发过程中可以没有编译这个环节，相较于 C 语言这门中级程序设计语言而言，Python 是一种高级程序设计语言，编程者运用起来更加容易

理解和便捷。

- Python 是交互式语言

这意味着，开发者在开发环境中写入代码，即刻就可以得到回馈的结果。在这个过程中，程序无须先进行整体的编译处理，而是逐条执行程序语句给出运行的结果。

- Python 是面向对象的程序设计语言

这意味着 Python 支持将代码封装成对象的面向对象的程序设计方式，使得 Python 程序能够充分发挥面向对象程序设计技术的长处。

- Python 是对初学者友好的语言

对初学者来说，Python 语言清晰简洁的编程方式，对语法的要求比其他语言更加宽松且具有丰富的扩展功能，便于初学编程者学习和掌握，即学习曲线更短。

1.1.2 有趣的 Python 程序

在深入学习 Python 语言之前，先来感受一下 Python 语言实现的一些有趣的小程序。

【范例程序 1-1】在屏幕上画一条蛇（Python 这单词的英文意思就是蟒蛇）

范例程序 1-1 的代码

```
import turtle          # 导入 turtle 包
turtle.setup(650,350,200,200)
turtle.penup()        # 提起笔移动，不绘图
turtle.fd(-250)       # 前进
turtle.pendown()      # 边移动边绘图
turtle.pensize(25)
turtle.pencolor("green")
turtle.seth(-40)
for I in range(3):
    turtle.circle(40,80)
    turtle.circle(-40,80)
turtle.circle(40,80/2)
turtle.fd(10)
turtle.circle(50,100)
turtle.circle(-25,130)
turtle.fd(20*6/5)
turtle.done()
```

执行以上代码，结果如图 1.1 所示。相信大多数读者通过运行这个范例程序代码，已经看到了这条可爱的小青蛇在屏幕上一点点画出来的样子，是不是很有趣呢？

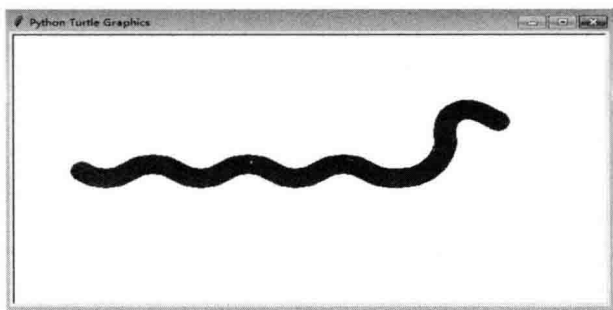


图 1.1 画蛇

同时,有些读者可能会对这个范例程序到底写了什么感到好奇。不用着急,请读者循序渐进、一点一点地弄明白这个范例程序代码中各种符号代表的含义。通过本书后面章节的学习,大家很快就可以自己轻松地画出想要的图形了。

Python 除了能够实现在屏幕上画图之外,还可以轻松解决一些数学问题。

【范例程序 1-2】实现简单的温度转换

目前有两种度量温度的标准,分别是摄氏度和华氏度,这个范例程序可以通过简单地输入一种温度值,然后经过换算得到另一种温度值。

范例程序 1-2 的代码

```
Temp = input("请输入带有符号的温度值: ")
if Temp[-1]in['F','f']:
    value = (eval(Temp[0:-1]) - 32 )/1.8
    print("转换后的温度值为{:.2f}C".format(value))
elif Temp[-1]in['C','c']:
    value = eval(Temp[0:-1])*1.8 +32
    print("转换后的温度值为{:.2f}F".format(value))
else:
    print("输入温度值格式有误")
```

尝试一下将这些程序代码复制到你的开发环境中,运行一下这个范例程序,试一试温度值的转换。范例程序 1-2 的运行结果如图 1.2 与图 1.3 所示。

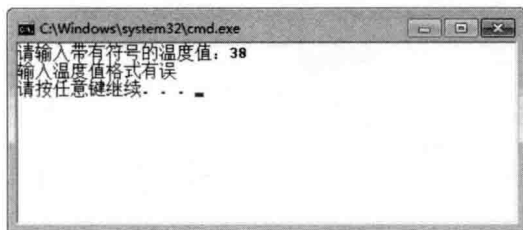


图 1.2 实现温度转换 I

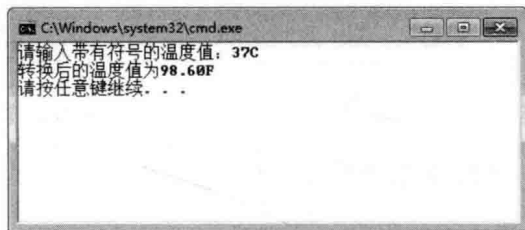


图 1.3 实现温度转换 II

有些读者可能还没有一个可以运行 Python 程序的集成开发环境,没关系,我们在下一节就指导读者安装和建立一个 Python 的集成开发环境。

1.2 集成开发环境

在上一节，大家看到了两个用 Python 语言编写的范例程序，心中一定充满了好奇，到底 Python 是怎么做到的呢？从本节开始，我们就正式开始 Python 的学习。俗话说，“工欲善其事，必先利其器。”所以在学习之前，我们需要了解一下 Python 的集成开发环境。

什么是集成开发环境？又有哪些可供选择的 Python 集成开发环境呢？

首先，集成开发环境就像画家手中的纸笔、厨师所用的炊具、木匠常用的工具，是开发程序必备的设施。设想一下，如果没了集成开发环境，那么程序要在哪里编写、执行和调试呢？

Python 的集成开发环境有很多，人们可以根据不同的需求来选择使用不同的集成开发环境。对于初学者来说，最适合的集成开发环境当然 Python 官网提供的。

1.2.1 安装 Python 3.7

下面先来安装 Python 3.7。

首先，读者需要知道自己所用的计算机配备的是什么操作系统。

Python 是一种跨平台的语言，可以在 Windows、Mac OS 以及各种 Linux 操作系统上运行。目前，人们在计算机中使用最多的还是微软公司的 Windows 操作系统，因此下面我们主要介绍如何在 Windows 上一步步安装 Python。

(1) 下载 Python 安装包。首先使用浏览器打开 Python 官网，网址为 <http://www.python.org/downloads/>，如图 1.4 所示。然后选择相应的安装包下载即可。

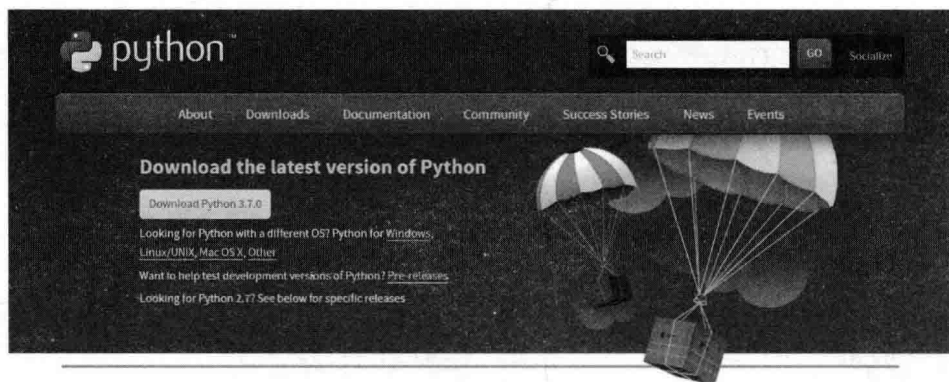


图 1.4 Python 官网

(2) 安装下载的 Python 3.7.2 安装包。执行安装包，出现类似图 1.5 所示的安装界面。

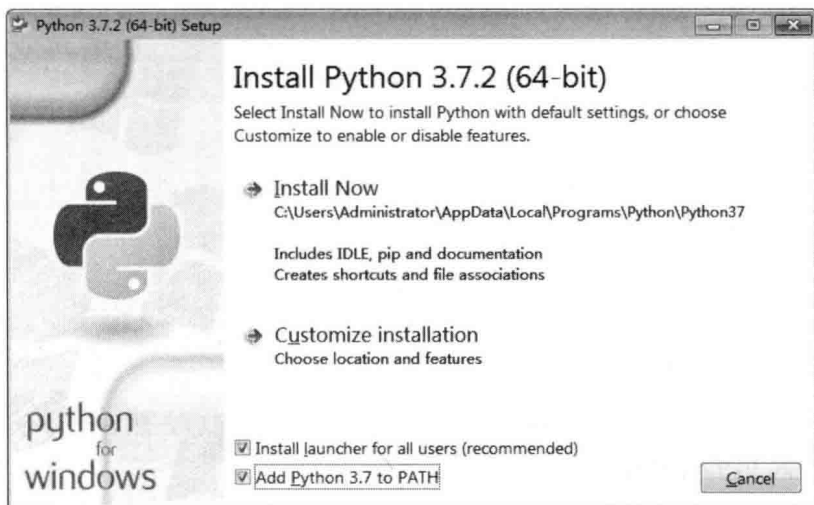


图 1.5 安装 Python 步骤 1

(3) 选择图 1.5 中的“Install Now”（开始安装）或者“Customize installation”（自定义安装）选项。在开始安装前需要先勾选“Add Python 3.7 to PATH”（将与运行 Python 3.7 有关的路径加入到系统环境变量 Path 中）选项。之后将出现如图 1.6 所示的安装界面。

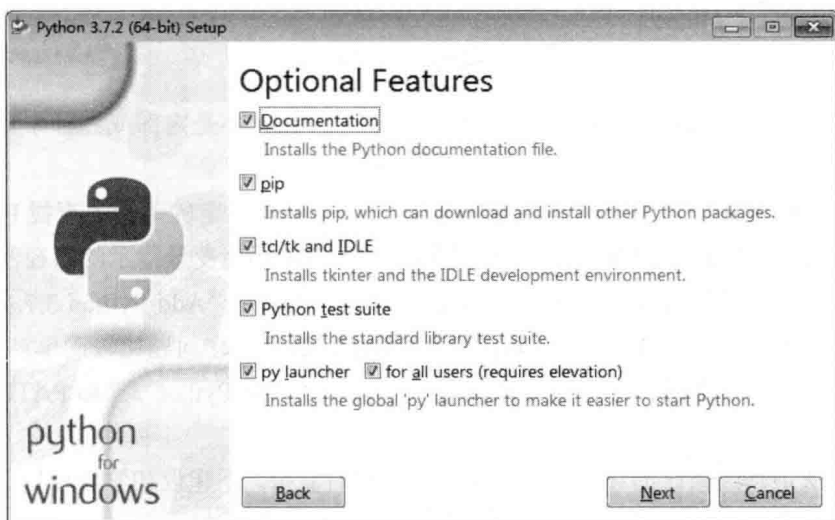


图 1.6 安装 Python 步骤 2

(4) 接下来单击【Next】按钮，直到完成。

(5) 测试安装是否成功。安装完毕之后需要对安装结果进行测试。首先打开系统的运行窗口，输入 cmd，如图 1.7 所示。



图 1.7 运行窗口

(6) 输入完成并单击【确定】按钮，之后会进入命令行窗口。然后在命令行窗口中输入“python”，进入 Python 开发环境，如图 1.8 所示。

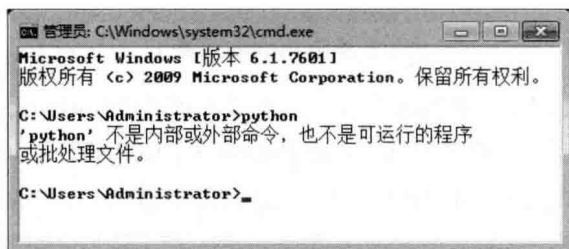


图 1.8 在命令行窗口中输入“python”

(7) 查看图 1.8，可以发现这里出现一个错误：‘python’不是内部或外部命令，也不是可运行的程序。

这是为什么呢？原因是 Windows 会根据安装时环境变量设定的路径去查找 Python 程序。出现图 1.8 所示的错误是因为 Windows 通过系统的 PATH 环境变量指示的路径没有找到这个程序。回忆一下，安装过程中是否漏掉了需要勾选的一个选项“Add Python 3.7 to PATH”。这时需要手动把安装的路径添加到 PATH 环境变量中。解决这个问题比较简单：

第一种是最简单的做法，即重新启动安装程序，将“Add Python 3.7 to PATH”选项勾选上。

第二种方法是修改系统的环境变量。这种方法相对复杂，不建议初学者使用。如果还是想尝试这种方法，建议读者去网络上查询修改 Path 环境变量的方法。

如果读者能够顺利地完成了上述安装步骤，那么恭喜读者已经获得了一把“趁手的武器”。接下来，就可以借助这个 Python 开发环境来编写 Python 程序了。

1.2.2 从 IDLE 启动 Python

启动 Python 系统自带的、简洁的集成开发环境 IDLE。

1. 通过运行 cmd 在命令行窗口中启动 IDLE

首先，找到 Windows 的“运行”选项，在运行框中输入“cmd”命令，启动命令行窗口，如图 1.9 所示。