

目录

第 1 章 入门准备

01 开篇词：你为什么要学 Python ？

02 我会怎样带你学 Python ？

03 让 Python 在你的电脑上安家落户

04 如何运行 Python 代码 ？

第 2 章 通用语言特性

05 数据的名字和种类—变量和类型

06 一串数据怎么存—列表和字符串

07 不只有一条路—分支和循环

08 将代码放进盒子—函数

09 知错能改—错误处理、异常机制

10 定制一个模子—类

11 更大的代码盒子—模块和包

12 练习—密码生成器

第 3 章 Python 进阶语言特性

13 这么多的数据结构（一）：列表、元祖、字符串

14 这么多的数据结构（二）：字典、集合

15 Python大法初体验：内置函数

16 深入理解下迭代器和生成器

17 生成器表达式和列表生成式

18 把盒子升级为豪宅：函数进阶

19 让你的模子更好用：类进阶

20 从小独栋升级为别墅区：函数式编程

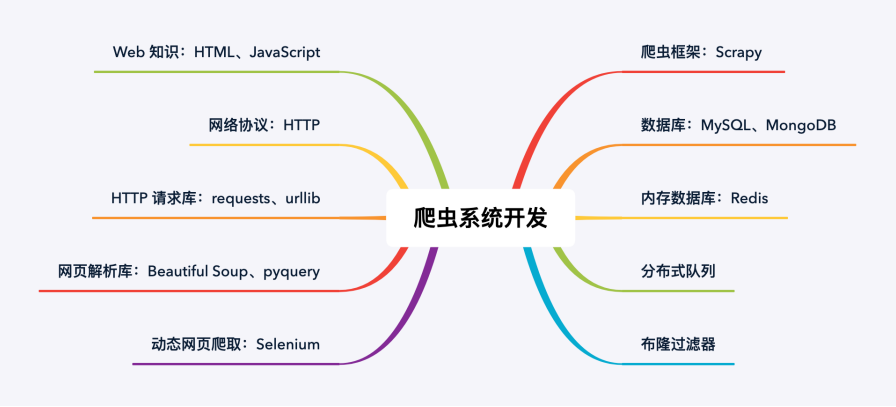
30 爬虫系统

更新时间：2019-11-06 09:33:26



“人不可有傲气，但不可无傲骨。”
——徐悲鸿

技术栈梳理



爬虫，即自动化的网页抓取程序，它能从网络中的大量网页里提取出所需的信息。网络是张大网，每个网页就是其中的一个交叉点，爬虫就顺着这张大网去爬取网页内容。

因为是从网页中提取内容，所以我们首先需要掌握基本的网页知识，如 HTML 和简单的 JavaScript。浏览器与 Web 服务器之间的通信采用 HTTP 协议，所以基本的 HTTP 知识也不能少。

爬虫需要模拟浏览器来向 Web 服务器发起请求，以获取网页内容。可以用 Python 的标准库 urllib，或者更好用的第三方库 requests 来达到这个目的。

拿到网页内容后，需要对网页进行解析，提取出其中的所需要的信息，或该网页上的其它网页链接。这时需要用到 Python 第三方库 BeautifulSoup 或 pyquery。

<div>← 慕课专栏</div> <div>≡ 你的第一本Python基础入门书 / 30 爬虫系统</div>	
目录	以上的抓取、解析过程也可以直接用专业的爬虫框架来完成，如 Scrapy，这是一种更工程化的方式。
第 1 章 入门准备	
01 开篇词：你为什么要学 Python？	当待抓取的网页数量很大时，这时需要注意网页判重的效率，也就是说需要有一种高效的方式来检查一个网页是否之前被抓取过，这就要用到布隆过滤器了。Redis 支持布隆过滤器扩展，能方便解决这个问题。
02 我会怎样带你学 Python？	
03 让 Python 在你的电脑上安家落户	当待抓取的网页数量进一步扩大时，单机的爬虫程序效率就十分低下了，需要考虑构建分布式的爬虫程序。也就是说在多台机器上同时来跑爬虫任务。我们需要一个分布式队列来统一管理、调度所有的抓取任务，Scrapy-Redis 可以做这件事。
04 如何运行 Python 代码？	
第 2 章 通用语言特性	入门资料
05 数据的名字和种类—变量和类型	了解和认识爬虫的最基本架构和运行流程，可以学习慕课网免费课程《Python开发简单爬虫》，或者阅读这篇知乎回答。
06 一串数据怎么存—列表和字符串	慕课网免费课程《Python最火爬虫框架Scrapy入门与实践》，可以带你了解和上手 Scrapy 的基本使用。Scrapy 的更详细内容可以阅读翻译文档《Scrapy入门教程》。
07 不只有一条路—分支和循环	MySQL 数据库入门，参看慕课网免费课程《与MySQL的零距离接触》。MongoDB 数据库入门，参看 W3Cschool 在线教程《MongoDB教程》。
08 将代码放进盒子—函数	至于前面所涉及技术的全面讲解，可以阅读图书《Python3网络爬虫开发实战》的免费在线版，基本上所有的爬虫相关知识和工具都有覆盖，值得一读。
09 知错能改—错误处理、异常机制	好了，资料不多，但要也要加油！
10 定制一个模子—类	
11 更大的代码盒子—模块和包	← 29 数据科学 31 机器学习 →
12 练习—密码生成器	
第 3 章 Python 进阶语言特性	精选留言 0
13 这么多的数据结构（一）：列表、元祖、字符串	欢迎在这里发表留言，作者筛选后可公开显示
14 这么多的数据结构（二）：字典、集合	
15 Python大法初体验：内置函数	
16 深入理解下迭代器和生成器	
17 生成器表达式和列表生成式	
18 把盒子升级为豪宅：函数进阶	
19 让你的模子更好用：类进阶	
20 从小独栋升级为别墅区：函数式编程	