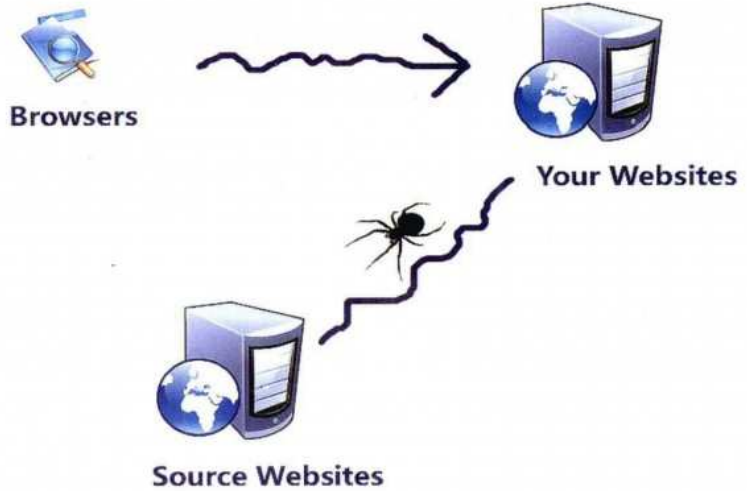
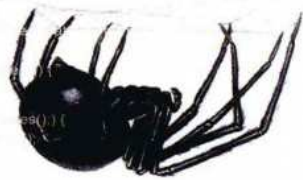


国内**第一本**专门讲解**网络爬虫**开发的书籍  
介绍如何应用**云计算**架构开发分布式爬虫



罗刚 王振东 编著

自己动手写

# 网络爬虫

- 猎兔搜索工程师多年项目经验总结
- 深入介绍Web数据挖掘实现过程
- 光盘中提供了高效的代码解决方案
- 案例均使用流行的Java语言编写



清华大学出版社

# 自己动手写 网络爬虫



本书介绍了网络爬虫开发中的关键问题与Java实现，主要包括从互联网获取信息、提取信息和对Web信息挖掘等内容，本书在介绍基本原理的同时注重通过具体代码实现来深入理解。书中部分代码甚至可以直接应用。本书适用于有Java程序设计基础的开发人员，同时也可以作为计算机相关专业本科生或研究生的参考教材。

ISBN 978-7-302-23647-4



9 787302 236474 >

定价：43.00元(附光盘1张)

# 前 言

当你在网上冲浪时，你是否知道还有一类特殊的网络用户也在互联网上默默地工作着，它们就是网络爬虫。这些网络爬虫按照设计者预定的方式，在网络中穿梭，同时自动收集有用的信息，进行分类和整理，将整理结果提供给用户，以方便用户查找他们感兴趣的内容。由于网络爬虫的实用性，引起了很多程序员，特别是 Web 程序员的兴趣。

但是大多数网络爬虫的开发原理与技巧在专业的公司内部都秘而不宣，至今仍然缺少理论与实践相结合的专门介绍网络爬虫的书籍。本书将弥补这个问题，尝试理论与实践相结合，深入透彻地讲解网络爬虫的原理，并且辅以相关代码作为参考。本书相关的代码在附带光盘中可以找到。

本书的两位主要作者在搜索引擎领域都有丰富的理论和实践经验。同时，还有多个程序员帮忙开发或编写了代码实现，例如 Java 实现异步 I/O 或对 PDF 文件的处理等。由于作者的日常工作繁忙，做得不够的地方敬请谅解。

作者罗刚在参加编写本书之前，还独立撰写过《自己动手写搜索引擎》一书，但存在讲解不够细致、知识点不够深入等问题。此次与王振东合著本书，相对于上一本书而言，对读者反馈有更高的预期。因为作者相信如下的假设：如果能够与更多的人更好地合作，事情往往能做得更好。

本书从基本的爬虫原理开始讲解，通过介绍优先级队列、宽度优先搜索等内容引领读者入门；之后根据当前风起云涌的云计算热潮，重点讲述了云计算的相关内容及其在爬虫中的应用，以及带偏好的爬虫、信息抽取、链接分析等内容；为了能够让读者更深入地了解爬虫，本书在最后两章还介绍了有关爬虫的数据挖掘的内容。

由于搜索引擎相关领域也正在快速发展中，而且由于篇幅的限制，有些不成熟的内容，没有能够在本书体现，例如有关“暗网”的内容。随着技术的不断发展，我们将在今后的版本中加入这些内容。

本书适合需要具体实现搜索引擎的程序员使用，对于信息检索等相关研究人员也有一定的参考价值，同时猎兔搜索技术团队也已经开发出以本书为基础的专门培训课程和商业软件。目前搜索引擎开发人员仍然很稀缺，作者真诚地希望通过本书把读者带入搜索引擎开发的大门并认识更多的朋友。

感谢开源软件和我们的家人、关心我们的老师和朋友、创业伙伴以及选择猎兔搜索软件的客户多年来的支持。

# 目 录

## 第 1 篇 自己动手抓取数据

第 1 章 全面剖析网络爬虫 .....	3	1.6 本章小结 .....	64
1.1 抓取网页 .....	4	第 2 章 分布式爬虫 .....	69
1.1.1 深入理解 URL .....	4	2.1 设计分布式爬虫 .....	70
1.1.2 通过指定的 URL 抓取 网页内容 .....	6	2.1.1 分布式与云计算 .....	70
1.1.3 Java 网页抓取示例 .....	8	2.1.2 分布式与云计算技术在爬虫 中的应用——浅析 Google 的 云计算架构 .....	71
1.1.4 处理 HTTP 状态码 .....	10	2.2 分布式存储 .....	72
1.2 宽度优先爬虫和带偏好的爬虫 .....	11	2.2.1 从 Ralation_DB 到 key/value 存储 .....	72
1.2.1 图的宽度优先遍历 .....	12	2.2.2 Consistent Hash 算法 .....	74
1.2.2 宽度优先遍历互联网 .....	13	2.2.3 Consistent Hash 代码实现 .....	79
1.2.3 Java 宽度优先爬虫示例 .....	15	2.3 Google 的成功之道——GFS .....	80
1.2.4 带偏好的爬虫 .....	22	2.3.1 GFS 详解 .....	80
1.2.5 Java 带偏好的爬虫示例 .....	23	2.3.2 开源 GFS——HDFS .....	84
1.3 设计爬虫队列 .....	24	2.4 Google 网页存储秘诀——BigTable .....	88
1.3.1 爬虫队列 .....	24	2.4.1 详解 BigTable .....	88
1.3.2 使用 Berkeley DB 构建 爬虫队列 .....	29	2.4.2 开源 BigTable——HBase .....	93
1.3.3 使用 Berkeley DB 构建爬虫 队列示例 .....	30	2.5 Google 的成功之道——MapReduce 算法 .....	98
1.3.4 使用布隆过滤器构建 Visited 表 .....	36	2.5.1 详解 MapReduce 算法 .....	100
1.3.5 详解 Heritrix 爬虫队列 .....	39	2.5.2 MapReduce 容错处理 .....	101
1.4 设计爬虫架构 .....	46	2.5.3 MapReduce 实现架构 .....	102
1.4.1 爬虫架构 .....	46	2.5.4 Hadoop 中的 MapReduce 简介 .....	104
1.4.2 设计并行爬虫架构 .....	47	2.5.5 wordCount 例子的实现 .....	105
1.4.3 详解 Heritrix 爬虫架构 .....	52	2.6 Nutch 中的分布式 .....	109
1.5 使用多线程技术提升爬虫性能 .....	55	2.6.1 Nutch 爬虫详解 .....	109
1.5.1 详解 Java 多线程 .....	55	2.6.2 Nutch 中的分布式 .....	116
1.5.2 爬虫中的多线程 .....	59	2.7 本章小结 .....	118
1.5.3 一个简单的多线程爬虫实现 ...	60		
1.5.4 详解 Heritrix 多线程结构 .....	61		



<b>第3章 爬虫的“方方面面”</b> .....	121	3.2.2 Java 主题爬虫 .....	128
3.1 爬虫中的“黑洞” .....	122	3.2.3 理解限定爬虫 .....	130
3.2 限定爬虫和主题爬虫 .....	122	3.2.4 Java 限定爬虫示例 .....	136
3.2.1 理解主题爬虫 .....	122	3.3 有“道德”的爬虫 .....	152
		3.4 本章小结 .....	155

## 第2篇 自己动手抽取 Web 内容

<b>第4章 “处理”HTML 页面</b> .....	159	5.3 抽取 RTF .....	218
4.1 征服正则表达式 .....	160	5.3.1 开源 RTF 文件解析器 .....	219
4.1.1 学习正则表达式 .....	160	5.3.2 实现一个 RTF 文件解析器 .....	219
4.1.2 Java 正则表达式 .....	164	5.3.3 解析 RTF 示例 .....	223
4.2 抽取 HTML 正文 .....	169	5.4 本章小结 .....	229
4.2.1 了解 HtmlParser .....	169	<b>第6章 多媒体抽取</b> .....	231
4.2.2 使用正则表达式抽取示例 .....	172	6.1 抽取视频 .....	232
4.3 抽取正文 .....	179	6.1.1 抽取视频关键帧 .....	232
4.4 从 JavaScript 中抽取信息 .....	194	6.1.2 Java 视频处理框架 .....	233
4.4.1 JavaScript 抽取方法 .....	195	6.1.3 Java 视频抽取示例 .....	237
4.4.2 JavaScript 抽取示例 .....	197	6.2 音频抽取 .....	249
4.5 本章小结 .....	199	6.2.1 抽取音频 .....	249
<b>第5章 非 HTML 正文抽取</b> .....	201	6.2.2 学习 Java 音频抽取技术 .....	253
5.1 抽取 PDF 文件 .....	202	6.3 本章小结 .....	256
5.1.1 学习 PDFBox .....	202	<b>第7章 去掉网页中的“噪声”</b> .....	257
5.1.2 使用 PDFBox 抽取示例 .....	206	7.1 “噪声”对网页的影响 .....	258
5.1.3 提取 PDF 文件标题 .....	207	7.2 利用“统计学”消除“噪声” .....	259
5.1.4 处理 PDF 格式的公文 .....	208	7.2.1 网站风格树 .....	262
5.2 抽取 Office 文档 .....	212	7.2.2 “统计学去噪”Java 实现 .....	270
5.2.1 学习 POI .....	212	7.3 利用“视觉”消除“噪声” .....	274
5.2.2 使用 POI 抽取 Word 示例 .....	213	7.3.1 “视觉”与“噪声” .....	274
5.2.3 使用 POI 抽取 PPT 示例 .....	215	7.3.2 “视觉去噪”Java 实现 .....	275
5.2.4 使用 POI 抽取 Excel 示例 .....	215	7.4 本章小结 .....	279

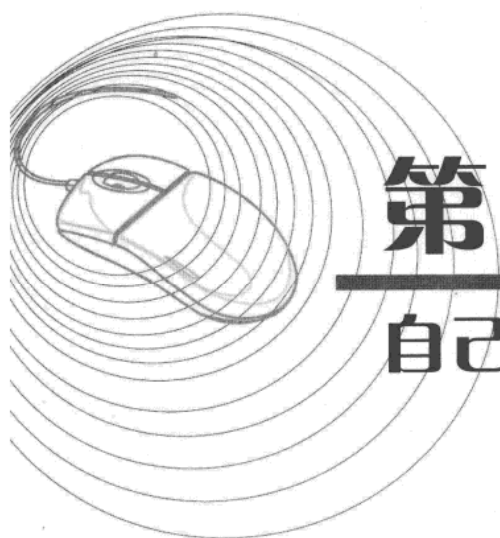
## 第3篇 自己动手挖掘 Web 数据

<b>第8章 分析 Web 图</b> .....	283	8.3.1 深入理解 PageRank 算法 .....	293
8.1 存储 Web “图” .....	284	8.3.2 PageRank 算法的 Java 实现 .....	297
8.2 利用 Web “图”分析链接 .....	293	8.3.3 应用 PageRank 进行	
8.3 Google 的秘密——PageRank .....	293	链接分析 .....	300



8.4 PageRank 的兄弟 HITS.....	301	9.6 本章小结 .....	331
8.4.1 深入理解 HITS 算法 .....	301	<b>第 10 章 分类与聚类的应用 .....</b>	<b>333</b>
8.4.2 HITS 算法的 Java 实现 .....	302	10.1 网页分类 .....	334
8.4.3 应用 HITS 进行链接分析 .....	313	10.1.1 收集语料库 .....	334
8.5 PageRank 与 HITS 的比较 .....	314	10.1.2 选取网页的“特征” .....	335
8.6 本章小结 .....	315	10.1.3 使用支持向量机进行 网页分类 .....	338
<b>第 9 章 去掉重复的“文档” .....</b>	<b>317</b>	10.1.4 利用 URL 地址进行 网页分类 .....	340
9.1 何为“重复”的文档 .....	318	10.1.5 使用 AdaBoost 进行 网页分类 .....	340
9.2 去除“重复”文档——排重 .....	318	10.2 网页聚类 .....	343
9.3 利用“语义指纹”排重 .....	318	10.2.1 深入理解 DBScan 算法 .....	343
9.3.1 理解“语义指纹” .....	320	10.2.2 使用 DBScan 算法 聚类实例 .....	344
9.3.2 “语义指纹”排重的 Java 实现 .....	321	10.3 本章小结 .....	346
9.4 SimHash 排重 .....	321		
9.4.1 理解 SimHash .....	322		
9.4.2 SimHash 排重的 Java 实现 .....	323		
9.5 分布式文档排重 .....	330		



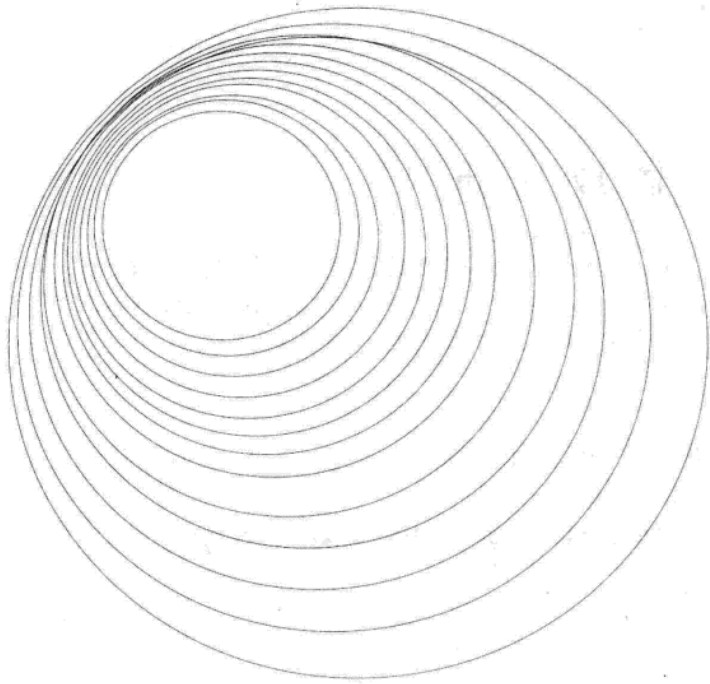


# 第 1 篇

---

## 自己动手抓取数据



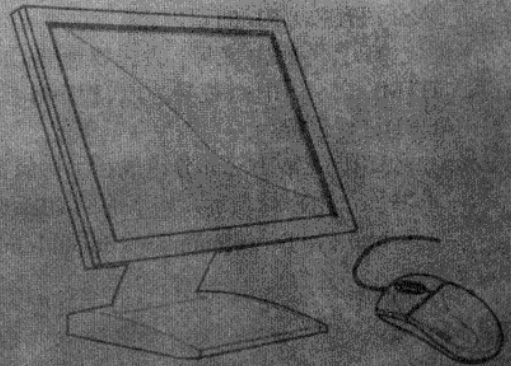


# 第 1 章

## 全面剖析网络爬虫

你知道百度、Google 是如何获取数以亿计的网页并且实时更新的吗？你知道在搜索引擎领域人们常说的 Spider 是什么吗？本章将全面介绍网络爬虫的方方面面。读完之后，你将完全有能力自己写一个网络爬虫，随意抓取互联网上任何感兴趣的东西。

既然百度、Google 这些搜索引擎巨头已经帮我们抓取了互联网上的大部分信息，为什么还要自己写爬虫呢？因为深入整合信息的需求是广泛存在的。在企业中，爬虫抓取下来的信息可以作为数据仓库多维展现的数据源，也可以作为数据挖掘的来源。甚至有人为了炒股，专门抓取股票信息。既然从美国中情局到普通老百姓都需要，那还等什么，让我们快开始吧。







## 1.1 抓取网页

网络爬虫的基本操作是抓取网页。那么如何才能随心所欲地获得自己想要的页面？这一节将从 URL 开始讲起，然后告诉大家如何抓取网页，并给出一个使用 Java 语言抓取网页的例子。最后，要讲一讲抓取过程中的一个重要问题：如何处理 HTTP 状态码。

### 1.1.1 深入理解 URL

抓取网页的过程其实和读者平时使用 IE 浏览器浏览网页的道理是一样的。比如，你打开一个浏览器，输入猎兔搜索网站的地址，如图 1.1 所示。

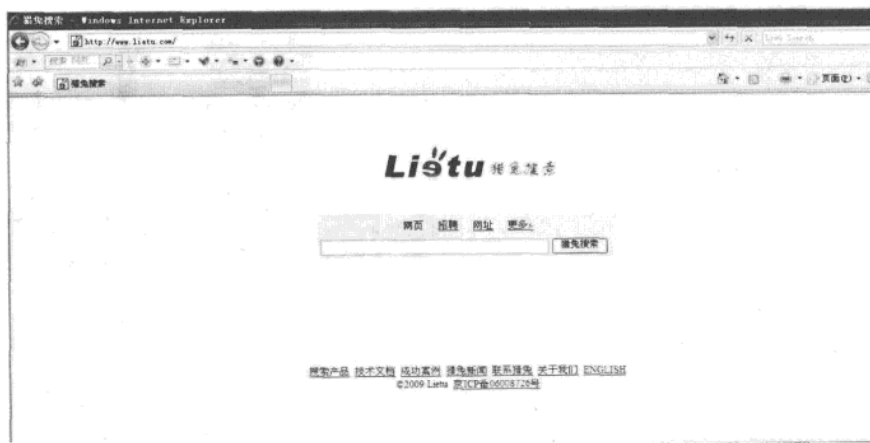


图 1.1 使用浏览器浏览网页

“打开”网页的过程其实就是浏览器作为一个浏览的“客户端”，向服务器端发送了一次请求，把服务器端的文件“抓”到本地，再进行解释、展现。更进一步，可以通过浏览器端查看“抓取”过来的文件源代码。选择“查看”|“源文件”命令，就会出现从服务器上“抓取”下来的文件的源代码，如图 1.2 所示。

在上面的例子中，我们在浏览器的地址栏中输入的字符串叫做 URL。那么，什么是 URL 呢？直观地讲，URL 就是在浏览器端输入的 `http://www.lietu.com` 这个字符串。下面我们深入介绍有关 URL 的知识。

在理解 URL 之前，首先要理解 URI 的概念。什么是 URI？Web 上每种可用的资源，如 HTML 文档、图像、视频片段、程序等都由一个通用资源标志符(Universal Resource Identifier, URI)进行定位。

URI 通常由三部分组成：①访问资源的命名机制；②存放资源的主机名；③资源自身的名称，由路径表示。如下面的 URI：

```
http://www.webmonkey.com.cn/html/html40/
```

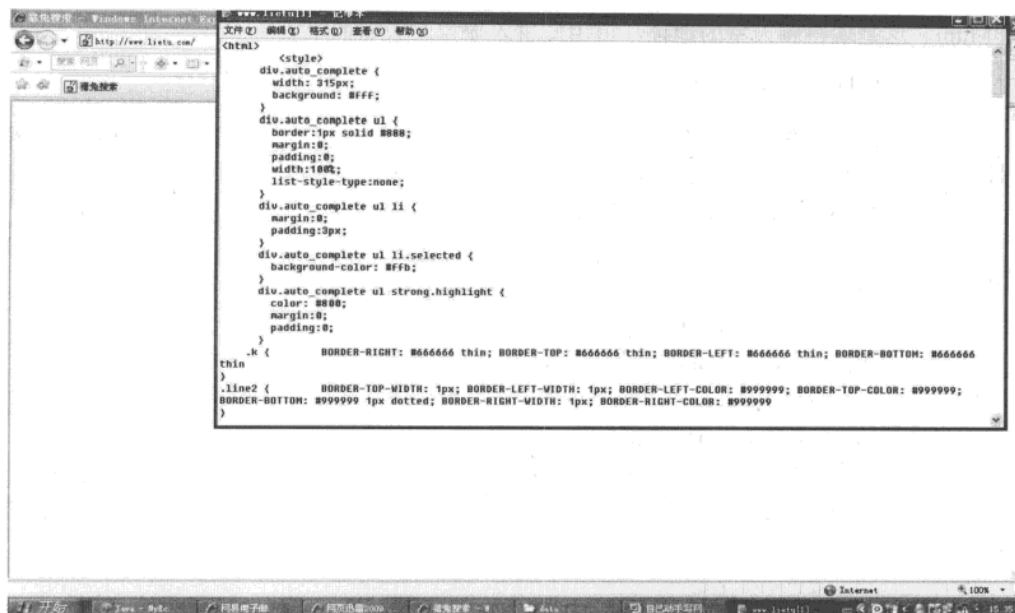


图 1.2 浏览器端源代码

我们可以这样解释它：这是一个可以通过 HTTP 协议访问的资源，位于主机 [www.webmonkey.com.cn](http://www.webmonkey.com.cn) 上，通过路径 “/html/html40” 访问。

URL 是 URI 的一个子集。它是 Uniform Resource Locator 的缩写，译为“统一资源定位符”。通俗地说，URL 是 Internet 上描述信息资源的字符串，主要用在各种 WWW 客户端程序和服务器程序上，特别是著名的 Mosaic。采用 URL 可以用一种统一的格式来描述各种信息资源，包括文件、服务器的地址和目录等。URL 的格式由三部分组成：

- 第一部分是协议(或称为服务方式)。
- 第二部分是存有该资源的主机 IP 地址(有时也包括端口号)。
- 第三部分是主机资源的具体地址，如目录和文件名等。

第一部分和第二部分用 “://” 符号隔开，第二部分和第三部分用 “/” 符号隔开。第一部分和第二部分是不可缺少的，第三部分有时可以省略。

根据 URL 的定义，我们给出了常用的两种 URL 协议的例子，供大家参考。

### 1. HTTP 协议的 URL 示例

使用超级文本传输协议 HTTP，提供超级文本信息服务的资源。

例：<http://www.peopledaily.com.cn/channel/welcome.htm>

其计算机域名为 [www.peopledaily.com.cn](http://www.peopledaily.com.cn)。超级文本文件(文件类型为 .html)是在目录 /channel 下的 [welcome.htm](http://www.peopledaily.com.cn/channel/welcome.htm)。这是中国人民日报的一台计算机。

例：<http://www.rol.cn.net/talk/talk1.htm>

其计算机域名为 [www.rol.cn.net](http://www.rol.cn.net)。超级文本文件(文件类型为 .html)是在目录 /talk 下的

欢迎访问：电子书学习和下载网站 (<https://www.shgis.com>)

文档名称：《自己动手写网络爬虫》罗刚，王振东 著.pdf

请登录 <https://shgis.com/post/4075.html> 下载完整文档。

手机端请扫码查看：

