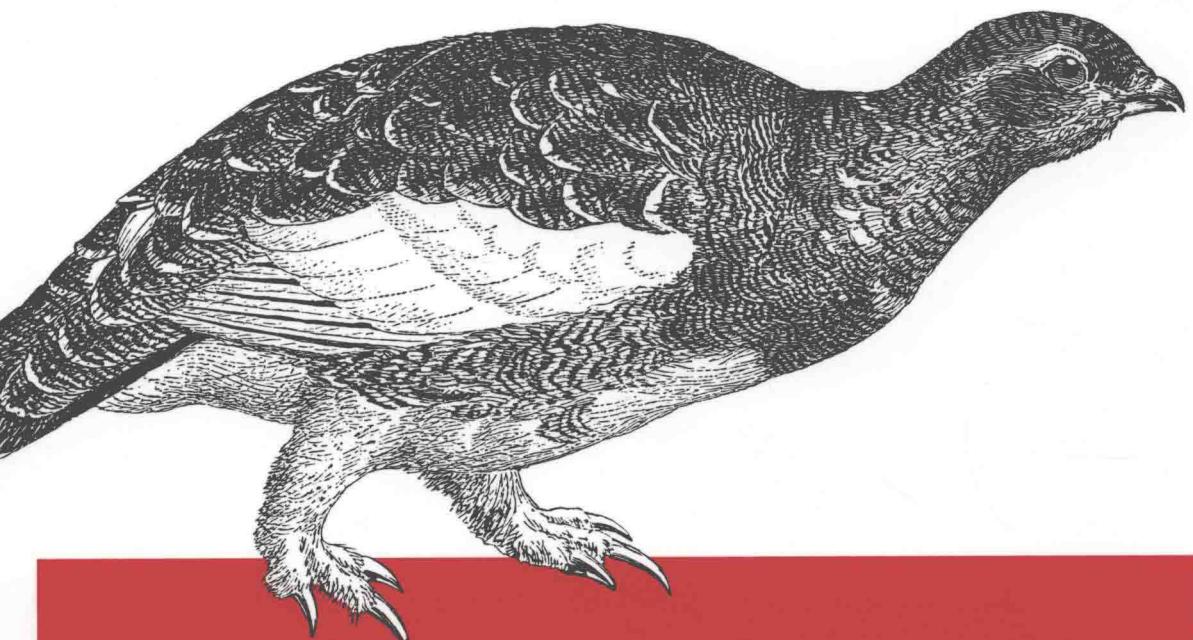


O'REILLY®

TURING

图灵程序设计丛书



# 数据科学入门

Data Science from Scratch

[美] Joel Grus 著  
高蓉 韩波 译



中国工信出版集团



人民邮电出版社  
POSTS & TELECOM PRESS



图灵程序设计丛书

# 数据科学入门

Data Science from Scratch  
First Principles with Python

[美] Joel Grus 著

高蓉 韩波 译



O'REILLY®

Beijing • Cambridge • Farnham • Köln • Sebastopol • Tokyo  
O'Reilly Media, Inc.授权人民邮电出版社出版

人民邮电出版社  
北京

## 图书在版编目（C I P）数据

数据科学入门 / (美) 格鲁斯 (Grus, J.) 著 ; 高蓉,  
韩波译. -- 北京 : 人民邮电出版社, 2016. 3  
(图灵程序设计丛书)  
ISBN 978-7-115-41741-1

I. ①数… II. ①格… ②高… ③韩… III. ①数据处  
理 IV. ①TP274

中国版本图书馆CIP数据核字(2016)第025033号

### 内 容 提 要

本书基于易于理解且具有数据科学相关的丰富的库的 Python 语言环境，从零开始讲解数据科学工作。具体内容包括：Python 速成，可视化数据，线性代数，统计，概率，假设与推断，梯度下降法，如何获取数据， $k$  近邻法，朴素贝叶斯算法，等等。作者借助大量具体例子以及数据挖掘、统计学、机器学习等领域的重要概念，详细展示了什么是数据科学。

本书适合有志成为数据科学工作者以及想了解数据科学的读者阅读。

- 
- ◆ 著 [美] Joel Grus
  - 译 高 蓉 韩 波
  - 责任编辑 朱 巍
  - 执行编辑 张 曼
  - 责任印制 杨林杰
  - ◆ 人民邮电出版社出版发行 北京市丰台区成寿寺路11号
  - 邮编 100164 电子邮件 315@ptpress.com.cn
  - 网址 <http://www.ptpress.com.cn>
  - 北京鑫正大印刷有限公司印刷
  - ◆ 开本：800×1000 1/16
  - 印张：19
  - 字数：451千字 2016年3月第1版
  - 印数：1-3 500册 2016年3月北京第1次印刷
  - 著作权合同登记号 图字：01-2015-8112号
- 

定价：69.00元

读者服务热线：(010)51095186转600 印装质量热线：(010)81055316

反盗版热线：(010)81055315

广告经营许可证：京东工商广字第 8052 号

---

# 版权声明

© 2015 by O'Reilly Media, Inc.

Simplified Chinese Edition, jointly published by O'Reilly Media, Inc. and Posts & Telecom Press, 2016. Authorized translation of the English edition, 2015 O'Reilly Media, Inc., the owner of all rights to publish and sell the same.

All rights reserved including the rights of reproduction in whole or in part in any form.

英文原版由 O'Reilly Media, Inc. 出版，2015。

简体中文版由人民邮电出版社出版，2016。英文原版的翻译得到 O'Reilly Media, Inc. 的授权。此简体中文版的出版和销售得到出版权和销售权的所有者——O'Reilly Media, Inc. 的许可。

版权所有，未得书面许可，本书的任何部分和全部不得以任何形式重制。



# O'Reilly Media, Inc.介绍

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。超级极客们正在开创着未来，而我们关注真正重要的技术趋势——通过放大那些“细微的信号”来刺激社会对新科技的应用。作为技术社区中活跃的参与者，O'Reilly 的发展充满了对创新的倡导、创造和发扬光大。

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极客和高瞻远瞩的商业领袖，共同描绘出开创新产业的革命性思想。作为技术人士获取信息的选择，O'Reilly 现在还将先锋专家的知识传递给普通的计算机用户。无论是通过书籍出版、在线服务或者面授课程，每一项 O'Reilly 的产品都反映了公司不可动摇的理念——信息是激发创新的力量。

## 业界评论

“O'Reilly Radar 博客有口皆碑。”

——*Wired*

“O'Reilly 凭借一系列（真希望当初我也想到了）非凡想法建立了数百万美元的业务。”

——*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、有前途、需要学习的主题。”

——*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，而且有几次都是一闪即逝的机会，尽管大路也不错。”

——*Linux Journal*

# 前言

## 数据科学

有人称数据科学家为“21世纪头号性感职业”(<https://hbr.org/2012/10/data-scientist-the-sexiest-job-of-the-21st-century/>)。虽说如此称呼有些夸张，但这个名称对数据科学的推崇却一点也没错，这是一个蓬勃发展、前途无限的行业。很多分析师都预言，未来十年会需要比现在多得多的数据科学工作者。

那么，什么是数据科学？唯有正确理解数据科学，才能培养出数据科学家。根据广受业界赞誉的文氏图 (<http://drewconway.com/zia/2013/3/26/the-data-science-venn-diagram>)，数据科学是以下几个方面的交叉：

- 黑客技能
- 数学和统计学知识
- 专业技能

我原本很想写一本能涵盖以上三个方面的书，但很快意识到仅关于专业技能的撰写就会耗费上万页笔墨，于是及时放弃转而专注于前两个方面。我的目标有两个：一是帮助读者掌握从事数据科学工作所必需的黑客技能；二是帮助读者熟悉数学和统计学，这是数据科学的核心。

对一本书来说，这两个愿望有点大了。学习黑客技能的最好方法就是钻研技术。通过阅读本书，你可以理解我钻研技术的方式，但相同的方式对你未必最适合；你可以理解我使用的一些工具，但相同的工具对你来说未必最顺手；你可以理解我如何解决数据问题，但相同的方式对你来说未必最有效。举例的目的和希望是启发你以自己的方式和方法完成工作。本书涵盖的所有代码和数据都可以从 GitHub 上下载。

同样，学习数学的最好方式就是研习数学。当然本书并不是一部数学著作，我们在本书中大半也不会“研习数学”，我想强调的是数学知识对从事数据科学工作至关重要。不理解

概率、统计、线性代数，就无法真正开始数据科学工作。在需要的地方，书中会引入数学方程式、数学直觉、数学公理，以及借以阐释大数学思想的卡通漫画。有我在，别怕！

总之，数据科学相当有趣（尤其和税务筹划或者煤矿开采等其他工作相比）。

## 从零开始

很多很多的数据科学库、框架、模块、工具箱可以有效地实现数据科学大部分常见的（和不常见的）算法与技术。如果你是一位数据科学家，就会非常熟悉 NumPy、scikit-learn、pandas 以及其他库。这些库对数据科学工作至关重要。如果还没有真正理解数据科学，运用这些库也是开始数据科学工作的好方式。

在本书中，我们从零开始着手数据科学工作。这意味着为了获得更好的理解，我们需要自己亲手构建工具和实现算法。我花费了很多心思选择注释良好、简洁易读的实现范例。在大部分情形下，所建立的工具意义清晰但实用性有限，它们对规模较小的示例数据集运转良好，但对“网络级别”的数据集就束手无策了。

在全书中，我会向读者指出相应的库，用以将相应技术运用于大规模数据集，但本书中我们不会使用它们。

对学习数据科学，一直有这样一种积极的争论，即什么样的语言环境最好？许多人认为是统计语言 R。（我们说，他们错了。）还有一些人认为是 Java 或者 Scala。而我认为，Python 才是最佳选择！

对于学习和从事数据科学工作，Python 具有几大优势：

- 免费；
- 编程相对简单（尤其是也易于理解）；
- 具有很多数据科学相关的库。

我不敢说 Python 是我最爱的编程语言，因为的确存在其他一些更舒适、设计更棒、编程更有乐趣的语言。但是，每当着手一个新的数据科学项目时，我最终使用的是 Python；每当需要快速构建某个有效程序的原型时，我使用的是 Python；每当需要用简洁易懂的方式表达数据科学概念时，我使用的还是 Python。于是，本书也采用 Python。

但是，教授 Python 不是本书的目的（尽管通过学习本书你会学到一些 Python 知识）。本书会用一章快速介绍 Python 的重要特征，这些特征与本书目的紧密相关。倘若读者没有 Python 基础（或编程基础），那需要再补充阅读一些关于 Python 的入门指导。

本书数据科学导论的其余部分采取了类似的书写方式，在必要或需要阐明时才深入细节，否则省略细节留给读者自己去挖掘（或者在维基百科上查阅）。

过去我曾培训过许多数据科学家。不是每个人都会努力变成改变世界的明星级数据忍者，但所有人都通过培训成为了更棒的数据科学家。我越来越相信，任何拥有一定数学基础和编程技术的人，只要再匹配一些基本材料就可以从事数据科学工作。必需品是好奇心、勤奋工作的态度，还有本书。没错，就是本书！

## 本书排版约定

本书使用了下列排版约定。

- 楷体  
表示新术语。
- 等宽字体 (**constant width**)  
表示程序片段，以及正文中出现的变量、函数名、数据库、数据类型、环境变量、语句和关键字等。
- 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**)  
表示应该由用户输入的值或根据上下文确定的值替换的文本。



该图标表示提示或建议。



该图标表示一般注释。



该图标表示警告或警示。

## 示例代码的使用

本书的补充材料（示例代码、练习等）都可以从 GitHub 下载：<https://github.com/joelgrus/>

data-science-from-scratch。

本书提供代码的目的是帮你快速完成工作。一般情况下，你可以在你的程序或文档中使用本书中的代码，而不必取得我们的许可，除非你想复制书中很大一部分代码。例如，你在编写程序时，用到了本书中的几个代码片段，这不必取得我们的许可。但若将 O'Reilly 图书中的代码制成光盘并进行出售或传播，则需获得我们的许可。引用示例代码或书中内容来解答问题无需许可。将书中很大一部分的示例代码用于你个人的产品文档，这需要我们的许可。

如果你引用了本书的内容并标明版权归属声明，我们对此表示感谢，但这不是强制的。版权归属声明通常包括：标题、作者、出版社和 ISBN，例如：“*Data Science from Scratch* by Joel Grus (O'Reilly). Copyright 2015 Joel Grus, 978-1-4919-0142-7”。

如果你认为你对示例代码的使用已经超出上述范围，或者你对是否需要获得示例代码的授权还不清楚，请随时联系我们：[permissions@oreilly.com](mailto:permissions@oreilly.com)。

## Safari® Books Online



Safari Books Online (<http://www.safaribooksonline.com>) 是应运而生的数字图书馆。它同时以图书和视频的形式出版世界顶级技术和商务作家的专业作品。技术专家、软件开发人员、Web 设计师、商务人士和创意专家等，在开展调研、解决问题、学习和认证培训时，都将 Safari Books Online 视作获取资料的首选渠道。

对于组织团体 (<https://www.safaribooksonline.com/enterprise/>)、政府机构 (<https://www.safaribooksonline.com/government/>)、教育机构 (<https://www.safaribooksonline.com/academic-public-library/>) 和个人，Safari Books Online 提供各种产品组合和灵活的定价策略 (<https://www.safaribooksonline.com/pricing/>)。用户可通过一个功能完备的数据库检索系统访问 O'Reilly Media、Prentice Hall Professional、Addison-Wesley Professional、Microsoft Press、Sams、Que、Peachpit Press、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 以及其他几十家出版社 (<https://www.safaribooksonline.com/our-library/>) 的上千种图书、培训视频和正式出版之前的书稿。要了解 Safari Books Online 的更多信息，我们网上见。

## 联系我们

请把对本书的评价和问题发给出版社。

美国：

O'Reilly Media, Inc.  
1005 Gravenstein Highway North  
Sebastopol, CA 95472

中国：

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/0636920033400.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](mailto: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

首先，我深深地感谢 Mike Loukides 接受我写作本书的提议（并对本书的篇幅提出了合理的建议）。其实他可以选择更轻松的方式，例如可以说：“那个总发样章给我的是什么人？我该怎样拒绝他？”但他没有，我感激不尽！同样，我很感谢我的编辑 Marie Beaugureau，她在整个出版过程中给予我指导，并最终使本书呈现出比我独立完成更棒的状态。

如果我从未学习数据科学，怎么能写出这样一本书？如果没有 Dave Hsu、Igor Tatarinov、John Rauser 和 Forecast 群组其他人的影响，我不大可能学习数据科学（当时甚至还没有数据科学这个名称）。还要感谢免费大型公开课程项目 Coursera 的优秀教师们。

同样，我很感谢本书的试读者和评论者。Jay Fundling 找出了许多错误，并指出很多含混的表达，因而本书得到了很大的改善，非常感谢。Debashis Ghosh 全面检查了本书的统计学部分。本书原本表达了肯定 Python 否定 R 的看法，Andrew Musselman 建议淡化这种立场，我后来体会到这是金玉良言。我也非常感谢 Trey Causey、Ryan Matthew Balfanz、Loris Mularoni、Núria Pujol、Rob Jefferson、Mary Pat Campbell、Zach Geary 和 Wendy Grus

给我的宝贵反馈。当然，本书其余的问题，责任在我。

我非常感谢 Twitter 上的数据科学社区，它让我接触到很多新奇的概念，让我认识了很多大牛。我深感自己的欠缺，需要写一本书来弥补。（再次）特别感谢 Trey Causey，他（并非刻意地）提醒我在书中加一章线性代数的内容。同样感谢 Sean J. Taylor，他（并非刻意地）指出了数据处理一章中的若干重大缺漏。

最后，向 Ganga 和 Madeline 致以我无限的感恩。比写一本书更痛苦的事，莫过于和写这本书的人朝夕相对。本书得以完成，全赖家人的支持与鼓励。

# 目录

前言	xiii
<b>第 1 章 导论</b>	1
1.1 数据的威力	1
1.2 什么是数据科学	1
1.3 激励假设：DataScencester	2
1.3.1 寻找关键联系人	3
1.3.2 你可能知道的数据科学家	5
1.3.3 工资与工作年限	8
1.3.4 付费账户	10
1.3.5 兴趣主题	11
1.4 展望	12
<b>第 2 章 Python 速成</b>	13
2.1 基础内容	13
2.1.1 Python 获取	13
2.1.2 Python 之禅	14
2.1.3 空白形式	14
2.1.4 模块	15
2.1.5 算法	16
2.1.6 函数	16
2.1.7 字符串	17
2.1.8 异常	18
2.1.9 列表	18
2.1.10 元组	19

2.1.11	字典	20
2.1.12	集合	22
2.1.13	控制流	23
2.1.14	真和假	24
2.2	进阶内容	25
2.2.1	排序	25
2.2.2	列表解析	25
2.2.3	生成器和迭代器	26
2.2.4	随机性	27
2.2.5	正则表达式	28
2.2.6	面向对象的编程	28
2.2.7	函数式工具	29
2.2.8	枚举	31
2.2.9	压缩和参数拆分	31
2.2.10	args 和 kwargs	32
2.2.11	欢迎来到 DataSciencester	33
2.3	延伸学习	33
<b>第 3 章</b>	<b>可视化数据</b>	<b>34</b>
3.1	matplotlib	34
3.2	条形图	36
3.3	线图	40
3.4	散点图	41
3.5	延伸学习	44
<b>第 4 章</b>	<b>线性代数</b>	<b>45</b>
4.1	向量	45
4.2	矩阵	49
4.3	延伸学习	51
<b>第 5 章</b>	<b>统计学</b>	<b>53</b>
5.1	描述单个数据集	53
5.1.1	中心倾向	55
5.1.2	离散度	56
5.2	相关	58
5.3	辛普森悖论	60
5.4	相关系数其他注意事项	61
5.5	相关和因果	62
5.6	延伸学习	63

<b>第 6 章 概率</b>	64
6.1 不独立和独立	64
6.2 条件概率	65
6.3 贝叶斯定理	66
6.4 随机变量	68
6.5 连续分布	68
6.6 正态分布	69
6.7 中心极限定理	72
6.8 延伸学习	74
<b>第 7 章 假设与推断</b>	75
7.1 统计假设检验	75
7.2 案例：掷硬币	75
7.3 置信区间	79
7.4 P-hacking	80
7.5 案例：运行 A/B 测试	81
7.6 贝叶斯推断	82
7.7 延伸学习	85
<b>第 8 章 梯度下降</b>	86
8.1 梯度下降的思想	86
8.2 估算梯度	87
8.3 使用梯度	90
8.4 选择正确步长	90
8.5 综合	91
8.6 随机梯度下降法	92
8.7 延伸学习	93
<b>第 9 章 获取数据</b>	94
9.1 stdin 和 stdout	94
9.2 读取文件	96
9.2.1 文本文件基础	96
9.2.2 限制的文件	97
9.3 网络抓取	99
9.3.1 HTML 和解析方法	99
9.3.2 案例：关于数据的 O'Reilly 图书	101
9.4 使用 API	105
9.4.1 JSON (和 XML)	105
9.4.2 使用无验证的 API	106
9.4.3 寻找 API	107

9.5 案例：使用 Twitter API .....	108
9.6 延伸学习 .....	111
<b>第 10 章 数据工作 .....</b>	<b>112</b>
10.1 探索你的数据 .....	112
10.1.1 探索一维数据 .....	112
10.1.2 二维数据 .....	114
10.1.3 多维数据 .....	116
10.2 清理与修改 .....	117
10.3 数据处理 .....	119
10.4 数据调整 .....	122
10.5 降维 .....	123
10.6 延伸学习 .....	129
<b>第 11 章 机器学习 .....</b>	<b>130</b>
11.1 建模 .....	130
11.2 什么是机器学习 .....	131
11.3 过拟合和欠拟合 .....	131
11.4 正确性 .....	134
11.5 偏倚 - 方差权衡 .....	136
11.6 特征提取和选择 .....	137
11.7 延伸学习 .....	138
<b>第 12 章 <math>k</math> 近邻法 .....</b>	<b>139</b>
12.1 模型 .....	139
12.2 案例：最喜欢的编程语言 .....	141
12.3 维数灾难 .....	146
12.4 延伸学习 .....	151
<b>第 13 章 朴素贝叶斯算法 .....</b>	<b>152</b>
13.1 一个简易的垃圾邮件过滤器 .....	152
13.2 一个复杂的垃圾邮件过滤器 .....	153
13.3 算法的实现 .....	154
13.4 测试模型 .....	156
13.5 延伸学习 .....	158
<b>第 14 章 简单线性回归 .....</b>	<b>159</b>
14.1 模型 .....	159
14.2 利用梯度下降法 .....	162
14.3 最大似然估计 .....	162
14.4 延伸学习 .....	163

<b>第 15 章 多重回归分析</b>	164
15.1 模型	164
15.2 最小二乘模型的进一步假设	165
15.3 拟合模型	166
15.4 解释模型	167
15.5 拟合优度	167
15.6 题外话：Bootstrap	168
15.7 回归系数的标准误差	169
15.8 正则化	170
15.9 延伸学习	172
<b>第 16 章 逻辑回归</b>	173
16.1 问题	173
16.2 Logistic 函数	176
16.3 应用模型	178
16.4 拟合优度	179
16.5 支持向量机	180
16.6 延伸学习	184
<b>第 17 章 决策树</b>	185
17.1 什么是决策树	185
17.2 熵	187
17.3 分割之熵	189
17.4 创建决策树	190
17.5 综合运用	192
17.6 随机森林	194
17.7 延伸学习	195
<b>第 18 章 神经网络</b>	196
18.1 感知器	196
18.2 前馈神经网络	198
18.3 反向传播	201
18.4 实例：战胜 CAPTCHA	202
18.5 延伸学习	206
<b>第 19 章 聚类分析</b>	208
19.1 原理	208
19.2 模型	209
19.3 示例：聚会	210
19.4 选择聚类数目 $k$	213

19.5	示例：对色彩进行聚类.....	214
19.6	自下而上的分层聚类.....	216
19.7	延伸学习.....	221
<b>第 20 章 自然语言处理 .....</b>		<b>222</b>
20.1	词云.....	222
20.2	n-grams 模型 .....	224
20.3	语法.....	227
20.4	题外话：吉布斯采样.....	229
20.5	主题建模.....	231
20.6	延伸学习.....	236
<b>第 21 章 网络分析 .....</b>		<b>237</b>
21.1	中介中心度.....	237
21.2	特征向量中心度.....	242
21.2.1	矩阵乘法.....	242
21.2.2	中心度 .....	244
21.3	有向图与 PageRank.....	246
21.4	延伸学习.....	248
<b>第 22 章 推荐系统 .....</b>		<b>249</b>
22.1	手工甄筛.....	250
22.2	推荐流行事物.....	250
22.3	基于用户的协同过滤方法.....	251
22.4	基于物品的协同过滤算法.....	254
22.5	延伸学习.....	256
<b>第 23 章 数据库与 SQL .....</b>		<b>257</b>
23.1	CREATE TABLE 与 INSERT .....	257
23.2	UPDATE .....	259
23.3	DELETE .....	260
23.4	SELECT .....	260
23.5	GROUP BY .....	262
23.6	ORDER BY .....	264
23.7	JOIN .....	264
23.8	子查询 .....	267
23.9	索引 .....	267
23.10	查询优化 .....	268
23.11	NoSQL .....	268
23.12	延伸学习 .....	269