

Natural Language Processing: Python and NLTK

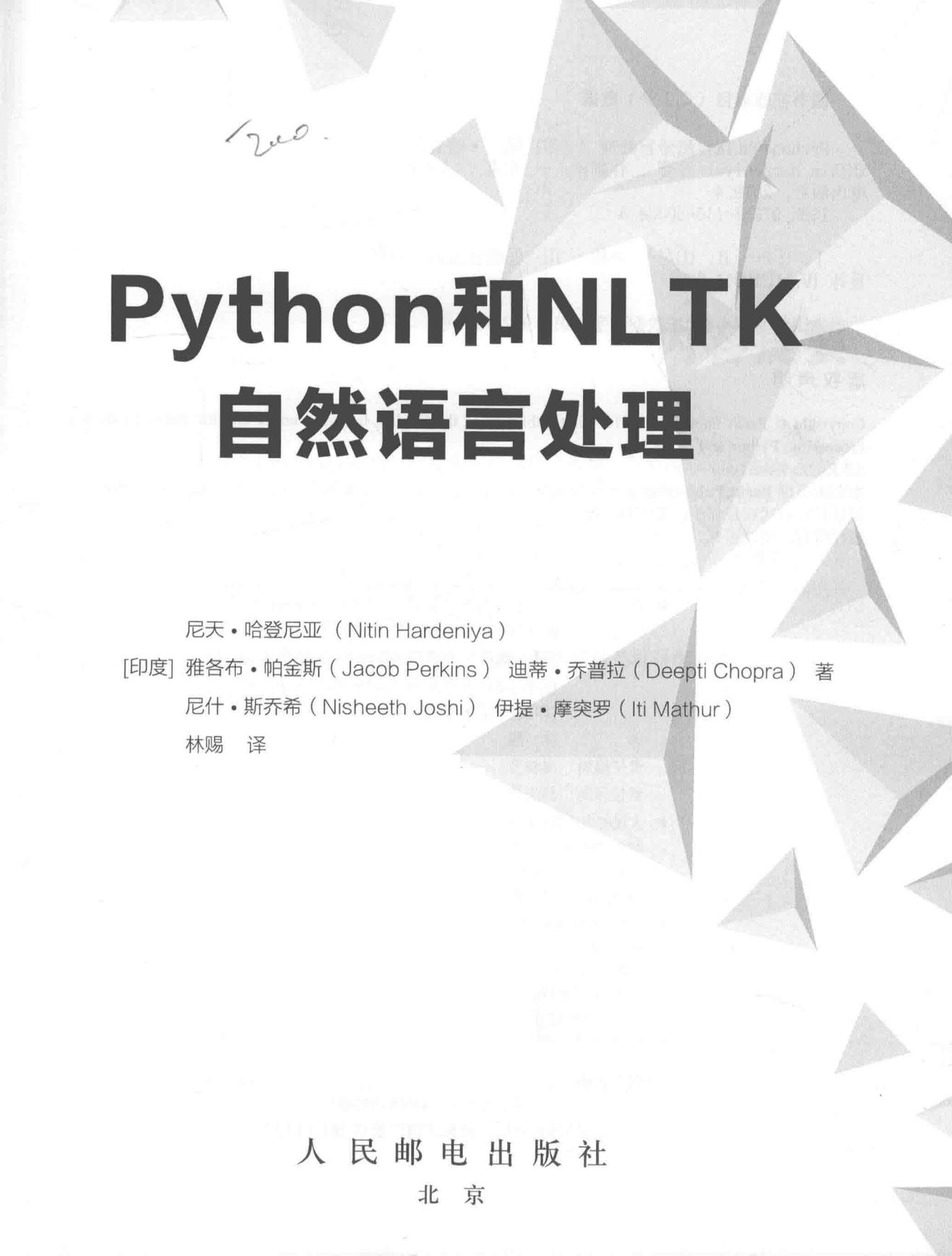
Python和NLTK 自然语言处理

尼天·哈登尼亚 (Nitin Hardeniya)

[印度] 雅各布·帕金斯 (Jacob Perkins) 迪蒂·乔普拉 (Deepti Chopra) 著

尼什·斯乔希 (Nisheeth Joshi) 伊提·摩突罗 (Iti Mathur)

林赐 译



1200.

Python和NLTK

自然语言处理

尼天·哈登尼亚 (Nitin Hardeniya)

[印度] 雅各布·帕金斯 (Jacob Perkins) 迪蒂·乔普拉 (Deepti Chopra) 著

尼什·斯乔希 (Nisheeth Joshi) 伊提·摩突罗 (Iti Mathur)

林赐 译

人民邮电出版社

北京

图书在版编目 (C I P) 数据

Python和NLTK自然语言处理 / (印) 尼天·哈登尼亚
(Nitin Hardeniya) 等著 ; 林赐译. — 北京 : 人民邮
电出版社, 2019. 4
ISBN 978-7-115-50334-3

I. ①P… II. ①尼… ②林… III. ①软件工具—程序
设计 IV. ①TP311.561

中国版本图书馆CIP数据核字(2018)第277916号

版权声明

Copyright © Packt Publishing 2016. First published in the English language under the title Natural Language Processing: Python and NLTK.

All Rights Reserved.

本书由英国 **Packt Publishing** 公司授权人民邮电出版社出版。未经出版者书面许可, 对本书的任何部分不得以任何方式或任何手段复制和传播。

版权所有, 侵权必究。

◆ 著 [印度]尼天·哈登尼亚 (Nitin Hardeniya)
[印度]雅各布·帕金斯 (Jacob Perkins)
[印度]迪蒂·乔普拉 (Deepti Chopra)
[印度]尼什·斯乔希 (Nisheeth Joshi)
[印度]伊提·摩突罗 (Iti Mathur)

译 林 赐
责任编辑 谢晓芳
责任印制 焦志炜

◆ 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>
固安县铭成印刷有限公司印刷

◆ 开本: 800×1000 1/16
印张: 40.75
字数: 810千字
印数: 1-2000册

2019年4月第1版

2019年4月河北第1次印刷

著作权合同登记号 图字: 01-2017-5038 号



定价: 138.00 元

读者服务热线: (010)81055410 印装质量热线: (010)81055316

反盗版热线: (010)81055315

广告经营许可证: 京东工商广登字 20170147 号

内容提要

本书旨在介绍如何通过 Python 和 NLTK 实现自然语言处理。本书包括三个模块。模块 1 介绍文本挖掘/NLP 任务中所需的所有预处理步骤，包括文本的整理和清洗、词性标注、对文本的结构进行语法分析、文本的分类等。模块 2 讲述如何使用 Python 3 的 NLTK 3 进行文本处理，包括标记文本、替换和校正单词、创建自定义语料库、词性标注、提取组块、文本分类等。模块 3 讨论了如何通过 Python 掌握自然语言处理，包括统计语言模型、词语形态学、词性标注、解析、语义分析、情感分析、信息检索等。

本书适合 NLP 和机器学习领域的爱好者、Python 程序员以及机器学习领域的研究人员阅读。

译者序

不下雪的日子，渥太华一片喧闹。2018年4月的雪刚刚化去，5月的郁金香节姗姗来迟，举目望去，大片大片的郁金香花把渥太华装点得五彩缤纷，绚烂多姿，生机盎然。在阳光的照射下，各色的郁金香花芬芳吐蕊、娇艳妩媚，令人陶醉和流连忘返。郁金香的芬芳还未散去，渥太华就迎来了“炎天暑月”的夏季，说是夏季，气温却很少超过30℃，此时，蛰伏了一个冬天的加拿大人，经受不住太阳的诱惑，蜂拥出动，三五成群地在里多运河的沙滩上玩耍嬉戏，好不自在。门户开放日（Doors Open Ottawa）、爵士音乐节、马拉松，各种活动精彩纷呈。正如鲁迅所说，越是民族的，就越是世界的。世界各地的移民带着自己家乡的文化和节日来到了渥太华。例如，中国的龙舟节成了这里的盛典，同时各个民族的人民也各自庆祝自己的节日，如夏至原住民节等。

加拿大是个移民国家，在这里，既可以看到不同的民族着装，也可以看到迥异的服饰风格。在不同的文化中，对个性的定义也各不相同，因此即使使用同一个词来形容一个人，也具有不同的韵味。作为中国人，我的思想不偏不倚，因此对各种奇装异服也就抱着听之任之的态度。我既欣赏热情奔放，也不排斥含蓄典雅。但是，对于计算机而言，要体察到语言的深层次含义，可谓“蜀道难，难于上青天”。本书介绍了如何编写能够理解人类语言的系统。

在翻译此书时，我想起了“通天塔”（The Tower of Babel）的故事。在《旧约》中，巴比伦人（Babylonian）想建造一座通天塔，但是上帝对此感到不满，因此借由变乱了人类的语言，使他们互相不能交流，破坏了这一工程。长久以来，人们一直在寻找不同语言之间的沟通方法。在20世纪三四十年代，凭借计算机的计算能力，人们开始了机器翻译的伟大尝试。这80多年来，人们在自然语言处理方面前赴后继，各个国家你追我赶，都想在这一领域拔得头筹，再造一座通天塔。直至最近，人类才第一次看到了梦想的曙光。市面上

也出现了各种机器翻译的软件，极大提高了人类译者的工作效率。这些软件具有可用性和市场价值。但是，机器翻译只是自然语言处理领域一个非常小的应用部分。自然语言处理其实包罗万象，让人眼花缭乱。其应用还包括问答系统、情感分析、图片题注、语音识别、词性标注、命名实体识别等。

自然语言处理博大精深，是一门融语言学、计算机科学、数学、统计于一体的科学。“纸上得来终觉浅，绝知此事要躬行。”要想完全掌握现在的自然语言技术，靠阅读几本书、几篇文章，是远远不够的。选择学习自然语言处理，并希望成为机器学习工程师，要经过三个阶段：实践、思考、创新。正如《荀子·儒效》中所说，“不闻不若闻之，闻之不若见之，见之不若知之，知之不若行之。学至于行止矣。”但是，“学而不思则罔，思而不学则殆。”在实践中，我们也要注意思考，寻找技术背后的原理，才能触类旁通，举一反三。作为科研工作者，仅仅成为知识的搬运工是不思进取的表现。“删繁就简三秋树，领异标新二月花。”掌握了知识之后，我们还应该学会创新，唯有这样，才可以成为社会的中流砥柱，引领技术的潮流。

在这里，我要特别感谢人民邮电出版社的领导和编辑，感谢他们对我的信任和理解，把这样一本好书交给我翻译。我也要感谢他们为本书的翻译投入了巨大的热情，可谓呕心沥血。没有他们的耐心和帮助，本书不可能顺利付梓。同时，在翻译过程中，加拿大友人 Jack Liu 和 Connie Wang 多次帮我指点迷津，才能使我为读者提供更贴切的译文。

译者才疏学浅，见闻浅薄，译稿中多有疏漏之处，还望读者谅解并不吝指正。读者如有任何意见和建议，请将反馈信息发送到邮箱 cilin2046@gmail.com。本书全由林赐翻译。

林 赐

前言

NLTK 是自然语言处理 (Natural Language Processing, NLP) 社区中最受欢迎和广泛使用的库之一。NLTK 的优点在于其简单性, 其中大多数复杂的 NLP 任务可以使用几行代码实现。本书主要内容包括: 如何将文本标记为各个单词, 如何使用 WordNet 语言词典, 如何以及何时进行词干提取或者词形还原, 如何替换单词和纠正拼写, 如何创建自己的自定义文本语料库和语料库 (包括 MongoDB 支持的语料库) 读取器, 如何使用词性标注器和部分词性标注单词, 如何使用部分解析创建和转换分块短语树, 如何进行文本分类的特征提取和情感分析, 如何进行并行和分布式文本处理, 以及如何在 Redis 中存储单词分布。

这种一边学习一边动手实践的学习方式会教你更多知识。本书有助于你成为使用 NLTK 进行自然语言处理的专家。

本书主要内容

模块 1 讨论文本挖掘/NLP 任务中所需的所有预处理步骤。该模块详细讨论标记化、词干提取、停用词删除和其他文本清理过程, 以及如何在 NLTK 中轻松实现这些操作。

模块 2 解释如何使用语料库读取器和创建自定义语料库。它还介绍如何使用 NLTK 附带的一些语料库。它涵盖组块过程 (也称为部分分析), 组块过程可以识别句子中的短语和命名实体。它还解释如何训练自己的自定义组块器并创建特定的命名实体识别器。

模块 3 讨论如何计算单词频率和实现各种语言建模技术。它还讨论浅层语义分析 (即 NER) 的概念和应用及使用 Wordnet 的 TSD。

模块 3 有助于你理解和应用信息检索与文本摘要的概念。

学习本书的软硬件配置

在学习模块 1 时，需要满足的软硬件配置如下表所示。

章号	需要的软件	免费/ 专用	下载软件的网站	硬件规格	需要的 操作系统
第 1~5 章	Python/Anaconda 和 NLTK	免费	Python 官网、continuum 官网 和 NLTK 官网	通用 UNIX 打印系统	任意
第 6 章	scikit-learn 和 gensim	免费	scikit-learn 官网、 radimrehurek 官网	通用 UNIX 打印系统	任意
第 7 章	Scrapy	免费	Scrapy 官网	通用 UNIX 打印系统	任意
第 8 章	NumPy、SciPy、 Pandas 和 Matplotlib	免费	Numpy 官网、Scipy 官网、 Pandas 官网和 Matplotlib 官网	通用 UNIX 打印系统	任意
第 9 章	Twitter、Python API 和 Facebook API	免费	Twitter 官网和 Facebook 官网	通用 UNIX 打印系统	任意

在学习模块 2 时，需要 Python 3 和列出的 Python 包。在本书中，作者使用了 Python 3.3.5。要安装这些包，可以使用 pip（参见 Python 官网）。以下是学习模块 2 时需要的包列表，以及编写本书时使用的版本号。

- NLTK $\geq 3.0a4$
- pyenchant $\geq 1.6.5$
- lock file $\geq 0.9.1$
- Numpy $\geq 1.8.0$
- Scipy $\geq 0.13.0$
- scikit-learn $\geq 0.14.1$
- execnet ≥ 1.1
- pymongo $\geq 2.6.3$

- Redis \geq 2.8.0
- lxml \geq 3.2.3
- BeautifulSoup4 \geq 4.3.2
- python-dateutil \geq 2.0
- Charade \geq 1.0.3

你还需要 NLTK-Trainer，可在 GitHub 网站上获得它。

除了 Python 之外，还有使用 MongoDB 和 Redis 两个 NoSQL 的技巧。MongoDB 可以在 MongoDB 官网下载，Redis 可以在 Redis 官网下载。

对于模块 3 的所有章节，使用了 Python 2.7 或 3.2+。NLTK 3.0 必须安装在 32 位计算机或 64 位计算机上。所需的操作系统是 Windows/Mac/UNIX 系统。

本书读者对象

如果你是 NLP 或机器学习爱好者以及想要快速掌握使用 NLTK 进行自然语言处理的中级 Python 程序员，那么本书的章节安排将为你带来很多好处。语言学和语义学分析方面的专业人士也会从本书中收益。

读者反馈

欢迎来自读者的反馈。让我们知道你对这本书的看法——你喜欢或不喜欢的内容。读者反馈对我们很重要，因为它可以帮助我们开发真正有用的选题。

要向我们发送一般反馈，请发送电子邮件至 feedback@packtpub.com，并在你的邮件标题中包含本书的书名。

如果你精通某方面的专业知识，并且有兴趣撰写或参与撰写某本书，请参考 packtpub.com 网站上的作者指南。

客户支持

既然你购买了 Packt 图书，那么还有更多配套资源可以帮助获得更大收益。

下载示例代码

可以从你在 packtpub.com 网站上的账号中下载本书的示例代码文件。如果你从其他地方购买了本书，那么请你访问 packtpub.com 网站并注册，之后客服人员会直接通过电子邮件向你发送示例代码文件。

可以按照以下步骤下载示例代码文件。

- (1) 使用你的电子邮件地址和密码登录 packtpub.com 网站并注册。
- (2) 将鼠标指针悬停在顶部的 SUPPORT 选项卡上。
- (3) 单击 Code Downloads & Errata 按钮。
- (4) 在 Search 框中输入书名。
- (5) 选择你要下载示例代码文件的书名。
- (6) 从你已购买书目的下拉菜单中选择对应的书名。
- (7) 单击 Code Download 按钮。

还可以通过单击 Packt Publishing 网站上本书页面上的 Code Files 按钮来下载代码文件。可以通过在 Search 框中输入本书的书名来访问此页面。请注意，你需要登录你的 Packt 账户。

下载文件后，请确保使用以下最新版本的解压缩软件解压缩文件夹。

- TinRAR/7-Zip（对于 Windows 系统）
- Zipeg/iZip/UnRarX（对于 Mac 系统）
- 7-Zip/PeaZip（对于 Linux 系统）

本书的代码包也托管在 GitHub 网站上。在 GitHub 网站上的书目、视频和课程目录中还提供了其他的代码包。请访问 GitHub 网站确认一下。

勘误表

虽然我们已尽力确保本书内容的准确性，但是错误在所难免。如果你在我们的一本书中发现了错误，可能是文字或代码中的错误，如果你能向提交勘误，我们将不胜感激。通过这样做，可以避免其他读者少走弯路，并帮助我们进一步提升本书后续版本的质量。

如果你发现了任何勘误，请访问 [packtpub](#) 官网，选择书名，单击 [Errata Submission Form](#) 链接，并输入详细的勘误信息进行报告。一旦你的勘误表通过了验证，将会接受你的提交，并且将勘误信息上传到我们的网站或添加到本书勘误部分下现有的勘误表中。

要查看以前提交的勘误表，请访问 [packtpub](#) 官网并在 Search 框中输入本书的名称。所需信息将显示在 Errata 部分下。

盗版行为

盗版因特网上受版权保护的内容是所有媒体上层出不穷的问题。在 Packt，我们非常重视保护我们的版权和许可。如果你在因特网发现以任何形式抄袭我们作品的行为，请立即向我们提供网络地址或网站名称，以便我们采取补救措施。

如果你发现了可疑的盗版内容，请过 copyright@packtpub.com 联系我们。

感谢你帮助保护作者的版权，能够为你提供有价值的内容，我们也感到非常欣慰。

问题

如果你对本书的任何方面有疑问，欢迎通过 questions@packtpub.com 联系我们，我们将尽力解决问题。

资源与支持

本书由异步社区出品，社区（<https://www.epubit.com/>）为您提供相关资源和后续服务。

配套资源

本书配套资源包括书中示例的源代码。

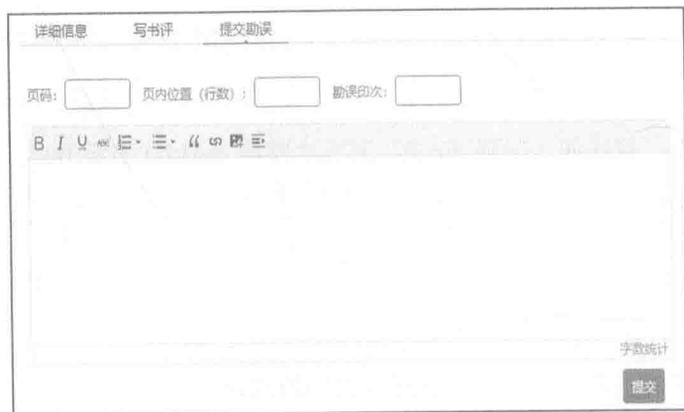
要获得以上配套资源，请在异步社区本书页面中单击 **配套资源**，跳转到下载界面，按提示进行操作即可。注意，为保证购书读者的权益，该操作会给出相关提示，要求输入提取码进行验证。

如果您是教师，希望获得教学配套资源，请在社区本书页面中直接联系本书的责任编辑。

提交勘误

作者和编辑尽最大努力来确保书中内容的准确性，但难免会存在疏漏。欢迎您将发现的问题反馈给我们，帮助我们提升图书的质量。

当您发现错误时，请登录异步社区，按书名搜索，进入本书页面，点击“提交勘误”，输入勘误信息，点击“提交”按钮即可。本书的作者和编辑会对您提交的勘误进行审核，确认并接受后，您将获赠异步社区的 100 积分。积分可用于在异步社区兑换优惠券、样书或奖品。



The screenshot shows a web form titled "提交勘误" (Submit勘误) with three tabs: "详细信息" (Detailed Information), "写书评" (Write a Review), and "提交勘误" (Submit勘误). The form contains three input fields: "页码:" (Page Number), "页内位置 (行数):" (Page Position (Line Number)), and "勘误印次:" (勘误印次). Below these fields is a rich text editor with a toolbar containing icons for bold (B), italic (I), underline (U), strikethrough (ABC), bulleted list, numbered list, link, unlink, and image. At the bottom right of the form, there is a "字数统计" (Character Count) label and a "提交" (Submit) button.

扫码关注本书

扫描下方二维码，您将会在异步社区微信服务号中看到本书信息及相关的服务提示。



与我们联系

我们的联系邮箱是 contact@epubit.com.cn。

如果您对本书有任何疑问或建议，请您发邮件给我们，请在邮件标题中注明本书书名，以便我们更高效地做出反馈。

如果您有兴趣出版图书、录制教学视频，或者参与图书翻译、技术审校等工作，可以发邮件给我们；有意出版图书的作者也可以到异步社区在线提交投稿（直接访问 www.epubit.com/selfpublish/submission 即可）。

如果您是学校、培训机构或企业，想批量购买本书或异步社区出版的其他图书，也可以发邮件给我们。

如果您在网上发现有针对异步社区出品图书的各种形式的盗版行为，包括对图书全部或部分内容的非授权传播，请您将怀疑有侵权行为的链接发邮件给我们。您的这一举动是对作者权益的保护，也是我们持续为您提供有价值的内容的动力之源。

关于异步社区和异步图书

“异步社区”是人民邮电出版社旗下 IT 专业图书社区，致力于出版精品 IT 技术图书和相关学习产品，为作译者提供优质出版服务。异步社区创办于 2015 年 8 月，提供大量精品 IT 技术图书和电子书，以及高品质技术文章和视频课程。更多详情请访问异步社区官网 <https://www.epubit.com>。

“异步图书”是由异步社区编辑团队策划出版的精品 IT 专业图书的品牌，依托于人民邮电出版社近 30 年的计算机图书出版积累和专业编辑团队，相关图书在封面上印有异步图书的 LOGO。异步图书的出版领域包括软件开发、大数据、AI、测试、前端、网络技术 etc。



异步社区



微信服务号

目录

模块 1 NLTK 基础知识

第 1 章 自然语言处理简介3

1.1 为什么要学习 NLP4

1.2 从 Python 的基本知识开始7

1.2.1 列表7

1.2.2 自助8

1.2.3 正则表达式9

1.2.4 词典11

1.2.5 编写函数11

1.3 NLTK13

1.4 试一试18

1.5 本章小结18

第 2 章 文本的整理和清洗19

2.1 文本整理19

2.2 文本清洗21

2.3 句子拆分离器22

2.4 标记解析22

2.5 词干提取24

2.6 词形还原25

2.7 停用词删除26

2.8 生僻字删除27

2.9 拼写校正27

2.10 试一试28

2.11 本章小结28

第 3 章 词性标注30

3.1 什么是词性标注30

3.1.1 斯坦福标注器33

3.1.2 深入了解标注器34

3.1.3 序列标注器35

3.1.4 布里尔标注器37

3.1.5 基于标注器的机器学习37

3.2 命名实体识别38

3.3 试一试40

3.4 本章小结41

第 4 章 对文本的结构进行语法分析42

4.1 浅层语法分析与深层语法
分析42

4.2 语法分析的两种方法43

4.3 为什么需要语法分析43

4.4 不同类型的语法分析器	45	6.2 文本分类	68
4.4.1 递归下降的语法分析器	45	6.3 采样	70
4.4.2 移位归约语法分析器	45	6.3.1 朴素贝叶斯	73
4.4.3 图表语法分析器	45	6.3.2 决策树	75
4.4.4 正则表达式语法 分析器	46	6.3.3 随机梯度下降	76
4.5 依存分析	47	6.3.4 逻辑回归	77
4.6 组块化	49	6.3.5 支持向量机	78
4.7 信息抽取	51	6.4 随机森林算法	79
4.7.1 命名实体识别	52	6.5 文本聚类	79
4.7.2 关系抽取	52	6.6 文本的主题建模	81
4.8 本章小结	53	6.7 参考资料	83
第5章 NLP 应用	54	6.8 本章小结	83
5.1 构建第一个 NLP 应用	54	第7章 网络爬取	85
5.2 其他的 NLP 应用	58	7.1 网络爬虫	85
5.2.1 机器翻译	58	7.2 编写第一个爬虫程序	86
5.2.2 统计机器翻译	59	7.3 Scrapy 中的数据流	89
5.2.3 信息检索	59	7.3.1 Scrapy 命令行界面	89
5.2.4 语音识别	61	7.3.2 项	94
5.2.5 文本分类	62	7.4 站点地图蜘蛛	96
5.2.6 信息提取	63	7.5 项管道	97
5.2.7 问答系统	64	7.6 外部参考	98
5.2.8 对话系统	64	7.7 本章小结	99
5.2.9 词义消歧	64	第8章 与其他 Python 库一同 使用 NLTK	100
5.2.10 主题建模	64	8.1 NumPy	100
5.2.11 语言检测	65	8.1.1 ndarray	101
5.2.12 光学字符识别	65	8.1.2 基本操作	102
5.3 本章小结	65	8.1.3 从数组中提取数据	103
第6章 文本分类	66	8.1.4 复杂的矩阵运算	103
6.1 机器学习	67	8.2 SciPy	107

8.2.1	线性代数	108	9.2	数据提取	126
8.2.2	特征值和特征向量	108	9.3	地理可视化	128
8.2.3	稀疏矩阵	109	9.3.1	影响者检测	129
8.2.4	优化	110	9.3.2	Facebook	130
8.3	Pandas	111	9.3.3	影响者的朋友	134
8.3.1	读取数据	112	9.4	本章小结	135
8.3.2	时序数据	114	第 10 章	大规模的文本挖掘	136
8.3.3	列转换	115	10.1	在 Hadoop 上使用 Python 的 不同方法	136
8.3.4	噪声数据	116	10.1.1	Python 的流	137
8.4	Matplotlib	117	10.1.2	Hive/Pig UDF	137
8.4.1	subplot	118	10.1.3	流包装器	137
8.4.2	添加轴	119	10.2	在 Hadoop 上运行 NLTK	138
8.4.3	散点图	120	10.2.1	UDF	138
8.4.4	柱状图	120	10.2.2	Python 流	140
8.4.5	3D 图	121	10.3	在 Hadoop 上运行 scikit-learn	141
8.5	外部参考	121	10.4	PySpark	144
8.6	本章小结	121	10.5	本章小结	146
第 9 章	使用 Python 进行社交媒体 挖掘	122			
9.1	数据收集	122			

模块 2 使用 Python 3 的 NLTK 3 进行文本处理

第 1 章	标记文本和 WordNet 的基础	149	1.3	将句子标记成单词	152
1.1	引言	149	1.3.1	工作方式	152
1.2	将文本标记成句子	150	1.3.2	工作原理	153
1.2.1	准备工作	150	1.3.3	更多信息	153
1.2.2	工作方式	151	1.3.4	请参阅	154
1.2.3	工作原理	151	1.4	使用正则表达式标记语句	154
1.2.4	更多信息	151	1.4.1	准备工作	155
1.2.5	请参阅	152	1.4.2	工作方式	155
			1.4.3	工作原理	155

1.4.4 更多信息	155	1.9.2 工作原理	165
1.4.5 请参阅	156	1.9.3 更多信息	166
1.5 训练语句标记生成器	156	1.9.4 请参阅	167
1.5.1 准备工作	156	1.10 发现单词搭配	167
1.5.2 工作方式	156	1.10.1 准备工作	167
1.5.3 工作原理	157	1.10.2 工作方式	167
1.5.4 更多信息	158	1.10.3 工作原理	168
1.5.5 请参阅	158	1.10.4 更多信息	168
1.6 在已标记的语句中过滤 停用词	158	1.10.5 请参阅	169
1.6.1 准备工作	158	第 2 章 替换和校正单词	170
1.6.2 工作方式	159	2.1 引言	170
1.6.3 工作原理	159	2.2 词干提取	170
1.6.4 更多信息	159	2.2.1 工作方式	171
1.6.5 请参阅	160	2.2.2 工作原理	171
1.7 查找 WordNet 中单词的 Synset	160	2.2.3 更多信息	171
1.7.1 准备工作	160	2.2.4 请参阅	173
1.7.2 工作方式	160	2.3 使用 WordNet 进行词形还原	173
1.7.3 工作原理	161	2.3.1 准备工作	173
1.7.4 更多信息	161	2.3.2 工作方式	173
1.7.5 请参阅	163	2.3.3 工作原理	174
1.8 在 WordNet 中查找词元和 同义词	163	2.3.4 更多信息	174
1.8.1 工作方式	163	2.3.5 请参阅	175
1.8.2 工作原理	163	2.4 基于匹配的正则表达式替换 单词	175
1.8.3 更多信息	163	2.4.1 准备工作	175
1.8.4 请参阅	165	2.4.2 工作方式	175
1.9 计算 WordNet 和 Synset 的 相似度	165	2.4.3 工作原理	176
1.9.1 工作方式	165	2.4.4 更多信息	177
		2.4.5 请参阅	177
		2.5 移除重复字符	177