



自然语言处理

理论与实战

唐聃 白宁超 冯暄 / 等著



中国工信出版集团



电子工业出版社
PUBLISHING HOUSE OF ELECTRONICS INDUSTRY
http://www.phei.com.cn

非外借

自然语言处理

理论与实战

唐聃 白宁超 冯暄 卿鸿宾 文俊 著

电子工业出版社

Publishing House of Electronics Industry

北京·BEIJING

内 容 简 介

自然语言处理是什么？谁需要学习自然语言处理？自然语言处理在哪些地方应用？相关问题一直困扰着不少初学者。针对这一情况，作者结合教学经验和工程应用编写此书。本书讲述自然语言处理相关学科知识和理论基础，并介绍使用这些知识的应用和工具，以及如何在实际环境中使用它们。由于自然语言处理的特殊性，其是一门多学科交叉的学科，初学者难以把握知识的广度和宽度，对侧重点不能全面掌握。本书针对以上情况，经过科学调研分析，选择以理论结合实例的方式将内容呈现出来。其中涉及开发工具、Python语言、线性代数、概率论、统计学、语言学等工程上常用的知识介绍，然后介绍自然语言处理的核心理论和案例解析，最后通过几个综合性的例子完成自然语言处理的学习和深入。本书旨在帮助读者快速、高效地学习自然语言处理和人工智能技术。

本书适用于具备一定编程基础的计算机专业、软件工程专业、通信专业、电子技术专业和自动化专业的大学二年级以上的学生、科研工作者和相关技术人员。一些做工程应用的自然语言处理工程师，也可以通过阅读本书补充理论知识，理论知识的魅力在于遇到工程难题时，可以知道其背后的原因，快速、准确地解决问题。

未经许可，不得以任何方式复制或抄袭本书之部分或全部内容。
版权所有，侵权必究。

图书在版编目（CIP）数据

自然语言处理理论与实战 / 唐聃等著. —北京：电子工业出版社，2018.6
ISBN 978-7-121-34390-2

I. ①自… II. ①唐… III. ①自然语言处理—研究 IV. ①TP391

中国版本图书馆 CIP 数据核字（2018）第 122905 号

责任编辑：陈晓猛

印 刷：三河市君旺印务有限公司

装 订：三河市君旺印务有限公司

出版发行：电子工业出版社

北京市海淀区万寿路 173 信箱

邮编：100036

开 本：787×980 1/16 印张：22.5

字数：432 千字

版 次：2018 年 6 月第 1 版

印 次：2018 年 6 月第 1 次印刷

定 价：79.00 元

凡所购买电子工业出版社图书有缺损问题，请向购买书店调换。若书店售缺，请与本社发行部联系，联系及邮购电话：（010）88254888，88258888。

质量投诉请发邮件至 zlts@phei.com.cn，盗版侵权举报请发邮件至 dbqq@phei.com.cn。

本书咨询联系方式：010-51260888-819，faq@phei.com.cn。

前言

本书讲述自然语言处理重要的相关学科知识和理论基础，并介绍使用这些知识的应用和工具，以及如何在实际环境中使用它们。市面上出版的自然语言处理书籍不多，且大多数讨论的是其背后的深奥原理，很少涉及基础知识和编程实现。自然语言处理是一门多学科交叉的学科，初入门的读者难以把握知识的广度和宽度，尤其对侧重点不能全面掌握。本书针对以上情况，经过科学调研分析，选择以理论结合实例的方式呈现知识点。首先介绍开发工具、Python 语言、线性代数、概率论、统计学、语言学等工程上常用的知识，然后介绍自然语言处理的核心理论和案例解析，最后通过几个综合性的例子完成自然语言处理的学习和深入。本书旨在帮助读者快速高效地学习自然语言处理和人工智能技术。

读者对象

自然语言处理是什么？谁需要学习自然语言处理？自然语言处理在哪些地方应用？本书就是对这几个问题的回答。自然语言处理领域主要探讨：如何处理及运用自然语言；自然语言认知（让计算机“懂”人类的语言）；自然语言生成系统（将计算机数据转化为自然语言）和自然语言理解系统（将自然语言转化为计算机程序更易于处理的形式）。自然语言处理在我们身边应用得非常广泛，其中包括：语音的自动合成与识别、机器翻译、自然语言理解、人机对话、信息检索、文本分类、自动文摘，等等。此外，自然语言处理也是人工智能、机器学习、深度学习的基础，重要程度不言而喻。如果读者有一定的编程基础，那么将有助于本书的阅读。如果读者不具备线性代数、概率论、统计学、语言学的知识，则可从本书中快速学习常见的工程应用知识；如果读者具备线性代数、概率论、统计学、语言学的知识，则更利于本书的阅读，可以对知识进行查全补充。此外，本身使用 Python 语言进行编程，假设读者具备 Python 知识，则可以跳过第 2 章，也更有利于本书的阅读。本书对于具备一定编程基础的计算机专业、软件工程专业、通信专业、电子技术专业和自动化专业的大学二年级以上的学生都是适宜的。一些做工程应用的自然语言处理工程师，也可以通过阅读本书补充理论知识。理论知识的最大魅力在于遇到工程难题时，可以知道其背后的原因，快速准确地解决问题。本书整体难度适宜，适合作为自学用书或课程教材。

本书结构

本书共四大部分 15 章，第一部分为基础部分，从第 1 章至第 6 章，主要介绍在自然语言交叉学科中，工程应用常用的学科知识，包括自然语言处理概述、Python 基础知识和环境搭建、线性代数、概率论、统计学、语言学。第二部分为理论部分，从第 7 章至第 14 章，主要介绍自然语言处理常用的理论知识，包括自然语言处理任务限制、技术范畴、语料库、中文自动分词、数据预处理、马尔可夫模型、条件随机场、模型评估和命名实体识别。第三部分为实战部分，第 15 章通过 GitHub 数据提取与可视化分析、微博话题爬取与存储分析，综合介绍网络爬虫、中文分词、数据处理、模型选择、数据分析、自然语言处理工具和数据可视化等技术点，这些技术也适用于以机器学习为代表的人工智能领域。本书各章节的具体内容介绍如下。

- ◎ 第 1 章基础入门：随着人工智能的快速发展，自然语言处理和机器学习技术的应用愈加广泛。然而身为初学者，要想快速入门这些前沿技术总是存在着各种各样的困难。为使读者对该领域的整体概况有一个系统明晰的认识，本章主要从发展历程、研究现状、应用前景等角度概要介绍自然语言处理及相关的机器学习技术。
- ◎ 第 2 章快速上手 Python：Python 作为一门简洁优美且功能强大的语言，越来越受到编程人员的青睐，在工业界和学术界也非常受欢迎。本书的全部代码都是通过 Python 实现的，之所以选择 Python 语言，是因为其可以跨平台跨应用开发，因此本章旨在帮助读者快速领略 Python 的概貌。如果读者已经具备 Python 基础，则可略过此章。
- ◎ 第 3 章线性代数：机器学习是计算机科学、统计学、数学和信息论等多个领域交叉的学科。线性代数又是数学的一个重要分支，对机器学习有着直接的影响。诸如算法建模、参数设置、验证策略、识别欠拟合和过拟合，等等。读者往往知道线性代数很有用，常常全书通读，造成时间不足和效率较低，归因于对线性代数在机器学习中的重点和用途不明。本章主要以简明的方式介绍常用的线性代数知识，并使读者知道线性代数常用于哪些方面。
- ◎ 第 4 章概率论：机器学习与深度学习是多学科交叉的科学技术，其中数学尤为重要，是很多形式化模型向数学建模的必经过程。继线性代数核心知识的介绍之后，本章着重介绍概率论的相关知识。
- ◎ 第 5 章统计学：在数据科学中，统计学的地位尤为显著。这是一门在数据分析的基础上，研究如何测定、收集、整理、归纳和分析数据规律，以便给出正确消息的学科。通过揭示数据背后的规律和隐藏信息，给相关角色提供参照价值，以做出相应的决策。其在数据挖掘、自然语言处理、机器学习中都被广泛应用。本章首先介绍常见的图形可视化的概念和使用，继而介绍数据度量标准、概率分布、统计假设检验、相关和回归，

以短小精悍的篇章使读者掌握基本的统计知识。

- ◎ 第 6 章语言学：本章主要从语音、词汇、语法三个角度对现代汉语进行一个简单概要的勾勒，在以往传统的语言学教材中一般还有“文字”“修辞”两节内容，因篇幅有限、与全书关联不强，在此删繁就简，未给读者一一呈现。需要注意的是，语言学本身是一门十分庞杂的学科，知识体系与研究方法或因语言不同而有区别，或因派别主义不同而有区别。但无论是何种语言，或是何门何派，在进行自然语言处理时我们要面临的永远是一个个真实的语料和具体的语言现象。理论是用来指导实践、拓宽我们研究思路的，究竟最后采用何种理论，这只是一个“白猫黑猫”的问题。
- ◎ 第 7 章自然语言处理：本章开篇直击要点，即自然语言处理的任务和限制。进而介绍其所涉及的主要技术范畴，并对这些技术方向进行介绍。在针对当前自然语言处理的难点进行详细剖析后，最终对 2017 年以后自然语言处理的发展进行展望。
- ◎ 第 8 章语料库：大数据发展的基石就是数据量的快速增加，无论是自然语言处理、数据挖掘、文本处理，还是机器学习领域，都是在此基础上通过规则或统计方法进行模型构建的。但是不是数据量足够大就叫大数据了呢？是不是数据量足够多就构成语料库了呢？带着这些疑问，本章将带你走进语料库的世界，对语料知识进行一次全面而深入的了解。
- ◎ 第 9 章中文自动分词：中文分词技术属于自然语言处理的技术范畴，中文分词是其他中文信息处理的基础，搜索引擎只是中文分词的一个应用。诸如机器翻译（MT）、语音合成、自动分类、自动摘要、自动校对，等等。
- ◎ 第 10 章数据预处理：数据预处理的整个步骤流程在自然语言处理的工程中要比其在机器学习的工程中精简一些，最大的区别就在于数据清洗和特征构造这两个至关重要的过程。在自然语言处理中特征构造是否良好，很大程度上取决于所构造的特征数据集的数据特性与文本内容语义吻合程度的高低。比如，文本情感分类和文本内容分类都属于分类范畴，但对于同一种算法（参数都调整到最优），在两个不同分类的业务下，得到的结果可能会相差很大。通过仔细分析，我们不难发现造成这种差异的根本原因就是构造出来的特征数据集的数据模式没有很好地契合文本的真实语义，这也是自然语言处理的最大难点。
- ◎ 第 11 章马尔可夫模型：笔者最早接触马尔可夫模型的定义源于吴军先生的《数学之美》一书，起初觉得深奥难懂且没什么用处。直到学习自然语言处理时，才真正使用到马尔可夫模型，并体会到此模型的奇妙之处。马尔可夫模型在处理序列分类时具有强大的功能，解决诸如词类标注、语音识别、句子切分、字素音位转换、局部句法剖

析、语块分析、命名实体识别、信息抽取等问题。此外它还广泛应用于自然科学、工程技术、生物科技、公用事业、信道编码等多个领域。

- ◎ 第 12 章条件随机场：条件随机场常用于序列标注、数据分割等自然语言处理任务中，此外在中文分词、中文人名识别和歧义消解等任务中也有应用。本书基于笔者在做语句识别序列标注过程中对条件随机场产生的了解。主要内容源于自然语言处理、机器学习、统计学习方法和部分网上资料对 CRF 的相关介绍，最后由笔者进行大量研究整理后汇总成知识体系。本章首先介绍条件随机场的相关概念，然后结合实例以期让读者深入理解条件随机场的应用。
- ◎ 第 13 章模型评估：本章源于基于 HMM 模型序列标注的一个实验，在实验完成之后，迫切想知道采用的序列标注模型好坏，有哪些指标可以度量。于是就产生了对这一专题进度的学习总结，这样也便于其他人参考。本章依旧简明扼要地梳理出模型评估核心指标，以期达到实用的目的。
- ◎ 第 14 章命名实体识别：命名实体识别在自然语言处理中占据着非常重要的地位，也是不可逾越的学术问题。命名实体识别的学术理论和研究方法众多，本章侧重整体介绍。首先阐述命名实体识别的背景知识和研究概况，介绍中文命名实体识别的特点与难点，辅以案例加深理解；然后对命名实体识别当前的研究方法和核心技术进行详细介绍；最后展望其在未来人工智能方面的发展前景。
- ◎ 第 15 章自然语言处理实战：自然语言处理技术是理论与实践相结合的一门学科，通过前面基础理论知识的介绍，读者对其理论有所认识，但其究竟有何用、怎么用却不深刻。本章通过实例演练，一方面对前面几章的知识进行复习回顾，另一方面利于加深理解研发的相关工作。本章的第一个案例以 GitHub 为例，实现数据提取和可视化；第二个案例以微博话题为例，实现数据采集、提取、存储与分析。

勘误

由于笔者能力有限，时间仓促，书中难免有错漏，欢迎读者批评指正。

联系方式：nlpjiaocheng@sina.com。

作者介绍

唐聃 教授，中科院工学博士。现工作于成都信息工程大学软件工程学院。研究方向包括自然语言处理、信息安全、数据分析。曾参与多项国家 863 项目和中科院知识创新工程项

目、省科技厅和教育厅项目；2016年入选中国科学院西部之光人才计划（中国科学院西部青年学者）。

白宁超 四川省计算机研究院软件开发工程师，曾参与多项四川省科技厅项目。其自然语言处理系列博文曾被CSDN、博客园、阿里云栖等多个技术社区转载。

冯暄 高级工程师，硕士学位。现任四川省计算机研究院信息化工程研究所所长。研究方向包括物联网、多源信息融合、软件工程。主持或参与国家级、省级科研项目16项。获得四川省科技进步奖二等奖2项、四川省科技进步奖三等奖1项。

卿鸿宾 四川大学中文系在校生。研究方向包括应用语言学、计算语言学、韵律句法学等。常年从事文学创作与文字工作，2017年作品《黄昏速写》发表于《子曰书院》微信公众号，取得了不错的反响。

文俊 硕士学位，现工作于成都广播电视台橙视传媒大数据中心，大数据算法工程师。研究方向主要包括数据挖掘、机器学习、自然语言处理、深度学习及云计算。

读者服务

轻松注册成为博文视点社区用户（www.broadview.com.cn），扫码直达本书页面。

- ◎ **提交勘误**：您对书中内容的修改意见可在 [提交勘误](#) 处提交，若被采纳，将获赠博文视点社区积分（在您购买电子书时，积分可用来抵扣相应金额）。
- ◎ **交流互动**：在页面下方 [读者评论](#) 处留下您的疑问或观点，与我们和其他读者一同学习交流。

页面入口：<http://www.broadview.com.cn/34390>



目录

第 1 章 基础入门	1
1.1 什么是自然语言处理.....	1
1.1.1 自然语言处理概述.....	1
1.1.2 自然语言处理的发展历史.....	3
1.1.3 自然语言处理的工作原理.....	6
1.1.4 自然语言处理的应用前景.....	7
1.2 开发工具与环境.....	7
1.2.1 Sublime Text 和 Anaconda 介绍.....	7
1.2.2 开发环境的安装与配置.....	8
1.3 实战：第一个小程序的诞生.....	13
1.3.1 实例介绍.....	13
1.3.2 源码实现.....	13
第 2 章 快速上手 Python	15
2.1 初识 Python 编程语言.....	15
2.1.1 Python 概述.....	15
2.1.2 Python 能做什么.....	17
2.1.3 Python 的语法和特点.....	19
2.2 Python 进阶.....	24
2.2.1 Hello World.....	24
2.2.2 语句和控制流.....	24
2.2.3 函数.....	27
2.2.4 List 列表.....	29
2.2.5 元组.....	32
2.2.6 set 集合.....	33
2.2.7 字典.....	33

2.2.8	面向对象编程：类.....	34
2.2.9	标准库.....	36
2.3	Python 深入——第三方库.....	36
2.3.1	Web 框架.....	36
2.3.2	科学计算.....	37
2.3.3	GUI.....	37
2.3.4	其他库.....	37
第 3 章	线性代数.....	39
3.1	线性代数介绍.....	39
3.2	向量.....	40
3.2.1	向量定义.....	40
3.2.2	向量表示.....	42
3.2.3	向量定理.....	42
3.2.4	向量运算.....	43
3.3	矩阵.....	47
3.3.1	矩阵定义.....	47
3.3.2	矩阵表示.....	48
3.3.3	矩阵运算.....	48
3.3.4	线性方程组.....	51
3.3.5	行列式.....	51
3.3.6	特征值和特征向量.....	55
3.4	距离计算.....	56
3.4.1	余弦距离.....	56
3.4.2	欧氏距离.....	57
3.4.3	曼哈顿距离.....	58
3.4.4	明可夫斯基距离.....	59
3.4.5	切比雪夫距离.....	61
3.4.6	杰卡德距离.....	62
3.4.7	汉明距离.....	63
3.4.8	标准化欧式距离.....	64
3.4.9	皮尔逊相关系数.....	65

第 4 章 概率论	67
4.1 概率论介绍	67
4.2 事件	68
4.2.1 随机试验	68
4.2.2 随机事件和样本空间	69
4.2.3 事件的计算	70
4.3 概率	71
4.4 概率公理	73
4.5 条件概率和全概率	76
4.5.1 条件概率	76
4.5.2 全概率	77
4.6 贝叶斯定理	78
4.7 信息论	79
4.7.1 信息论的基本概念	79
4.7.2 信息度量	80
第 5 章 统计学	85
5.1 图形可视化	85
5.1.1 饼图	85
5.1.2 条形图	88
5.1.3 热力图	91
5.1.4 折线图	93
5.1.5 箱线图	96
5.1.6 散点图	99
5.1.7 雷达图	102
5.1.8 仪表盘	104
5.1.9 可视化图表用法	106

5.2 数据度量标准	108
5.2.1 平均值	108
5.2.2 中位数	108
5.2.3 众数	110
5.2.4 期望	111
5.2.5 方差	112
5.2.6 标准差	113
5.2.7 标准分	114
5.3 概率分布	115
5.3.1 几何分布	115
5.3.2 二项分布	116
5.3.3 正态分布	118
5.3.4 泊松分布	121
5.4 统计假设检验	123
5.5 相关和回归	125
5.5.1 相关	125
5.5.2 回归	127
5.5.3 相关和回归的联系	130
第 6 章 语言学	132
6.1 语音	132
6.1.1 什么是语音	132
6.1.2 语音的三大属性	133
6.1.3 语音单位	134
6.1.4 记音符号	135
6.1.5 共时语流音变	136
6.2 词汇	137
6.2.1 什么是词汇	137
6.2.2 词汇单位	137
6.2.3 词的构造	138

6.2.4	词义及其分类	140
6.2.5	义项与义素	141
6.2.6	语义场	142
6.2.7	词汇的构成	143
6.3	语法	143
6.3.1	什么是语法	143
6.3.2	词类	144
6.3.3	短语	148
6.3.4	单句	150
6.3.5	复句	152
第 7 章	自然语言处理	155
7.1	自然语言处理的任务和限制	155
7.2	自然语言处理的主要技术范畴	156
7.2.1	语音合成	156
7.2.2	语音识别	156
7.2.3	中文自动分词	157
7.2.4	词性标注	158
7.2.5	句法分析	158
7.2.6	文本分类	159
7.2.7	文本挖掘	160
7.2.8	信息抽取	161
7.2.9	问答系统	161
7.2.10	机器翻译	162
7.2.11	文本情感分析	163
7.2.12	自动摘要	164
7.2.13	文字蕴涵	165

7.3 自然语言处理的难点.....	165
7.3.1 语言环境复杂.....	165
7.3.2 文本结构形式多样.....	166
7.3.3 边界识别限制.....	166
7.3.4 词义消歧.....	167
7.3.5 指代消解.....	168
7.4 自然语言处理展望.....	169
第 8 章 语料库.....	173
8.1 语料库浅谈.....	173
8.2 语料库深入.....	174
8.3 自然语言处理工具包：NLTK.....	176
8.3.1 NLTK 简介.....	176
8.3.2 安装 NLTK.....	177
8.3.3 使用 NLTK.....	180
8.3.4 在 Python NLTK 下使用 Stanford NLP.....	186
8.4 获取语料库.....	194
8.4.1 国内外著名语料库.....	195
8.4.2 网络数据获取.....	197
8.4.3 NLTK 获取语料库.....	200
8.5 综合案例：走进大秦帝国.....	208
8.5.1 数据采集和预处理.....	208
8.5.2 构建本地语料库.....	208
8.5.3 大秦帝国语料操作.....	209
第 9 章 中文自动分词.....	216
9.1 中文分词简介.....	216
9.2 中文分词的特点和难点.....	218
9.3 常见中文分词方法.....	219
9.4 典型中文分词工具.....	220
9.4.1 HanLP 中文分词.....	220

9.4.2 其他中文分词工具.....	223
9.5 结巴中文分词	224
9.5.1 基于 Python 的结巴中文分词.....	224
9.5.2 结巴分词工具详解.....	227
9.5.3 结巴分词核心内容.....	230
9.5.4 结巴分词基本用法.....	233
第 10 章 数据预处理	241
10.1 数据清洗.....	241
10.2 分词处理.....	242
10.3 特征构造.....	242
10.4 特征降维与选择.....	243
10.4.1 特征降维.....	243
10.4.2 特征选择.....	243
10.5 简单实例.....	244
10.6 本章小结.....	249
第 11 章 马尔可夫模型	250
11.1 马尔可夫链	250
11.1.1 马尔可夫简介	250
11.1.2 马尔可夫链的基本概念	251
11.2 隐马尔可夫模型.....	253
11.2.1 形式化描述	253
11.2.2 数学形式描述	255
11.3 向前算法解决 HMM 似然度.....	256
11.3.1 向前算法定义	256
11.3.2 向前算法原理	256
11.3.3 现实应用：预测成都天气的冷热	258
11.4 文本序列标注案例：Viterbi 算法	259

第 12 章 条件随机场	263
12.1 条件随机场介绍	263
12.2 简单易懂的条件随机场	265
12.2.1 CRF 的形式化表示	265
12.2.2 CRF 的公式化表示	266
12.2.3 深度理解条件随机场	268
第 13 章 模型评估	269
13.1 从统计角度介绍模型概念	269
13.1.1 算法模型	269
13.1.2 模型评估和模型选择	270
13.1.3 过拟合与欠拟合的模型选择	272
13.2 模型评估与选择	275
13.2.1 模型评估的概念	275
13.2.2 模型评估的评测指标	275
13.2.3 以词性标注为例分析模型评估	276
13.2.4 模型评估的几种方法	278
13.3 ROC 曲线比较学习器模型	279
第 14 章 命名实体识别	281
14.1 命名实体识别概述	281
14.2 命名实体识别的特点与难点	284
14.3 命名实体识别方法	284
14.4 中文命名实体识别的核心技术	286
14.5 展望	295
第 15 章 自然语言处理实战	296
15.1 GitHub 数据提取与可视化分析	296
15.1.1 了解 GitHub 的 API	296
15.1.2 使用 NetworkX 作图	299
15.1.3 使用 NetworkX 构建兴趣图	301

15.1.4	NetWorkX 部分统计指标.....	304
15.1.5	构建 GitHub 的兴趣图	305
15.1.6	可视化.....	318
15.2	微博话题爬取与存储分析	320
15.2.1	数据采集.....	320
15.2.2	数据提取.....	329
15.2.3	数据存储.....	332
15.2.4	项目运行与分析.....	333
附录 A	Python 与其他语言调用.....	337
附录 B	Git 项目上传简易教程.....	339
参考文献	341