

全方位解读大数据背景下的文本分析核心算法

助力打造下一代AI

# 在线文本 数据挖掘

算法原理与编程实现

刘通◎著



**要点**涵盖5大基础领域、6大应用场景

120个**模型**展示语义分析，技术场景化、具象化

精练150多个**图表**，快速掌握核心原理

通过60多个经典**案例**，高效引导商业实践

基于百万级语料分析**数据**，实战分析全面精准



中国工信出版集团

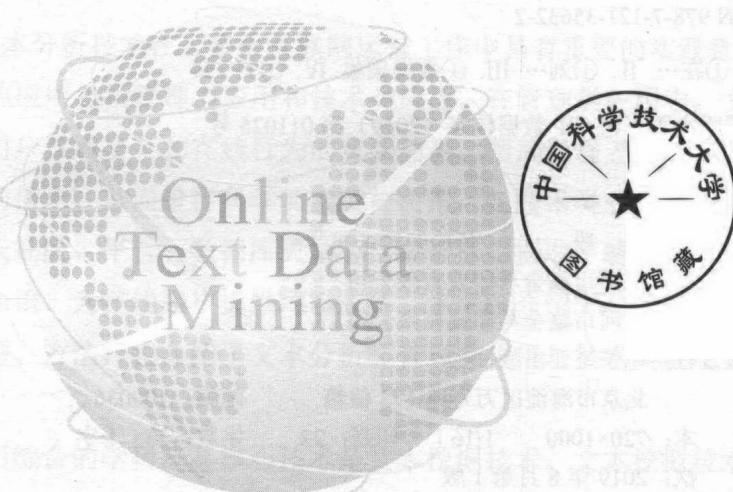


电子工业出版社  
PUBLISHING HOUSE OF ELECTRONICS INDUSTRY  
<http://www.phei.com.cn>

# 在线文本 数据挖掘

算法原理与编程实现

刘通◎著



电子工业出版社

Publishing House of Electronics Industry  
北京•BEIJING

## 内 容 简 介

为了满足大数据环境下网络运营与管理的需求，本书详细而系统地介绍了有关文本分析的核心技术与方法。本书基于统计分析、数据挖掘、机器学习等计算机技术，介绍了如何对在线环境的文本内容进行建模与分析，同时介绍了文本分析技术的具体应用场景。本书并非是纯粹的技术类书籍，而是一本教授读者如何更好地应用技术的实践手册。

本书分为 13 章，内容主要包括 3 个方面：①文本分析概要，包括概述、预备知识；②文本分析的基础类方法，包括文本建模、文本分类、文本聚类、序列标注；③文本分析的应用类方法，包括信息检索、文本摘要、口碑分析、社交网络分析、深度学习与 NLP、实证研究。

本书内容丰富、详略得当，结构清晰、系统。阅读本书需要读者具备一定的统计学知识和与数据挖掘相关的基础知识。本书特别适合对文本分析技术感兴趣的学生、科研工作者，以及数据分析类职业的工作人员阅读和参考。

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

在线文本数据挖掘：算法原理与编程实现 / 刘通著. —北京：电子工业出版社，2019.8

ISBN 978-7-121-35632-2

I . ①在… II . ①刘… III . ①数据采集 IV . ①TP274

中国版本图书馆 CIP 数据核字（2019）第 011025 号

责任编辑：张毅      特约编辑：田学清

印 刷：三河市鑫金马印装有限公司

装 订：三河市鑫金马印装有限公司

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：720×1000      1/16      印张：22      字数：388 千字

版 次：2019 年 8 月第 1 版

印 次：2019 年 8 月第 1 次印刷

定 价：88.00 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) 57565805。

## 作者简介

---



### 刘通

上海交通大学安泰经济与管理学院  
管理科学与工程专业博士，现任国内  
知名金融机构算法工程师。作者曾在  
管理学领域和计算机科学领域的国内  
核心期刊上发表过诸多与文本挖掘应  
用技术相关的重要研究成果。此外，  
作者在就读博士期间，将技术理论有  
效结合实践，采用文本挖掘方法对用  
户在线医疗平台上的医生选择和医疗  
在线服务购买的决策行为进行了深入  
研究，对如何改进在线平台的盈利模  
型并提高市场收益提出有价值的管理  
学见解。作者还曾参与过上海交通大  
学与华为技术有限公司合作的分布式  
计算平台的深度学习算法选型研究，  
比较了包括Theano、Tensorflow、  
Deeplearning4j在内的不同深度学习  
平台在分布式环境下的性能与算法兼  
容情况。



为读者出好书，传播文字的思想价值；

美迪出版

为作者做好书，提升智慧的商品价值。

团购电话：010-88254062

预知新书信息、交流投稿、邮购团购

请发邮件至：[xyan@phei.com.cn](mailto:xyan@phei.com.cn)

本书在全国各大新华书店、书城均有销售

浏览请登录：[www.phei.com.cn](http://www.phei.com.cn)

新浪微博@美迪出版

豆瓣：美迪出版

试读结束：需要全本请在线购买：[www.ertongbook.com](http://www.ertongbook.com)

# 前言

在大数据时代，数据的价值开始被推上各行各业的舞台。人们更注重从海量的数据中挖掘感兴趣的信息，以实现丰富的技术应用，进行科学的管理决策。在互联网环境中，数据的分析与利用尤为重要，尤其是数值类型数据的分析和文本类型数据的分析。其中，文本类型数据的分析比一般数值类型数据的分析复杂，文本类型数据是大数据 4V 特征的具体体现，其相关技术也更具难度。尽管如此，文本类型数据在整个网络中的信息占比仍十分庞大，且对用户的各种在线交互、活动及购买行为也有着不容小觑的影响。因此，网络中的文本类型数据具有十分重要的分析价值。本书将重点对当今文本类型数据的重要分析技术进行详细、系统的介绍。

在应用方面，文本分析技术在大多数互联网运营工作中具有重要的实践意义。基于文本分析技术的应用包括管理类应用和技术类应用。在管理类应用中，文本分析可以有效提取用户在线交互和在线行为的重要信息，帮助管理者更好地掌握用户、产品、市场的信息，从而进行科学的建模与决策；在技术类应用中，文本分析可以充分从在线社区、平台、数据库大量的文本数据中提取、解析、创造用户感兴趣的信息与知识，为在线用户提供内容服务。本书既介绍了与文本分析密切相关的理论、模型、方法，也介绍了文本分析在管理类应用与技术类应用等具体场景中的实现。

文本分析是一门综合的学科，其核心技术是文本挖掘技术。文本挖掘技术与传统的数据挖掘技术一脉相承，是数据挖掘在语言学领域中的应用。从事文本分析的数据分析者不仅需要掌握丰富的数据处理、建模及挖掘方法，还需要掌握语言学知识、社会学知识，也需要充分理解语言产生的背景、应用和使用语言信息的用户对象。文本数据比一般的数值数据更容易体现人类的感情与行为，其相应的技术也具备更高的智能化程度，因此，在任何领域，掌握文本分析技术对数据分析者来说都是一个不小的挑战。

近些年，随着整个信息社会对文本数据重视程度的提升，以及计算机软硬件技术的飞速发展，文本分析领域的研究成果形成井喷式爆发。由于篇幅所限，本书虽然无法全面讲解文本分析的所有前沿技术，但是仍然尽可能地将所有经典的、有代表性的研究成果展现给大家，使从事文本分析的工作人员、科研人员及文本分析技术的爱好者能够高效而系统地对整个文本分析领域有一定的了解。阅读本书后，希望读者能够具备基于文本分析技术的能力，从而解决工作中的各种文本分析问题，并能深刻地认识到文本分析为互联网领域及整个社会带来的实践价值。

## 本书特色

### 1. 内容丰富，系统全面，详略得当

本书内容涵盖了当前大部分主流的文本分析技术与方法，笔者按照自身的知识体系对其进行细致的归纳与梳理，并由浅入深地向读者进行了系统的介绍。本书内容详略得当，突出了知识的重点、难点。书中内容依托于数据技术，但不拘泥于技术本身，在介绍相关技术理论时注重向读者教授核心方法及思维方式，帮助读者掌握技术的核心理念，从而使读者做到灵活应用、深入思考、举一反三、即时实践。

### 2. 行文通俗易懂，随意而不失严谨，有利于读者快速吸收理解

本书在介绍知识时，尽可能地用通俗易懂的语言对技术细节进行描述，而不是生硬地对学术文献中的定义、规范和公式进行搬运。对于很多技术难点，笔者均赋予了自身的思考和感悟，并用生动而接地气的语言进行了转述。

本书中所有方法和理论都具有翔实可靠的学术依据，是科学而严谨的，所介绍的方法和技术也都得到了学术上的广泛认可和接受。本书还在特定的位置附注了关键知识点的学术来源，以供感兴趣的读者进一步进行知识的补充、考证。

### 3. 图文并茂，配备实例，有趣生动

本书虽是一本技术类书籍，但在排版风格上力争做到图文并茂，以增加读者的阅读兴趣，提高读者对于知识的理解效率。一图胜千字，本书中很多文本分析中重要的技术流程采用了示意图的表述方式，这可以有效地对知识点进行串联与总结。

此外，对于很多分析方法，本书还介绍了其具体应用场景，以及具体技术实现。这样，读者不仅掌握了知识的核心理念，根据具体实例也知道了如何运用知识。本书在知识结构上，可大致分成基础篇和应用篇，基础篇重点讲述理论方法，而应用篇偏向于知识在具体场景中的技术实现。本书在知识点设计方面更加生动灵活，有效地保证了文本分析技术的落地与推广。

## 本书内容及体系结构

### 第1章 概述

本章详细谈论了大数据时代下互联网公司的机会与挑战，介绍了在线文本分析技术在网站运营中重要的战略性地位。本章还基于大数据背景，从4V角度介绍了文本分析的主要技术特征。本章内容可以帮助读者更好地了解在线文本分析总体的知识框架和体系。

### 第2章 预备知识

本章引入了与在线文本分析密切相关的理论知识。首先，介绍了文本挖掘的主要任务，并介绍了与其相关的一些重要理论知识，如文本语义分析与语法分析、文本的结构化分析与标准化分析。其次，介绍了机器学习的基本概念，阐述了机器学习与深度学习的关系。对于机器学习，本章涉及的技术要点主要包括概率图模型、判别式模型、产生式模型、机器学习模型求解，以及模型过拟合。

### 第3章 文本建模

本章介绍了文本分析的基本任务——文本建模，即科学而有效地将非结构化的文本类型数据转换为可以直接进行数据分析与挖掘的数值类型数据。本章介绍了文本建模的主要应用场景，并从语言学建模和统计学建模两个主要方面对相关技术进行了详细介绍。

### 第4章 文本分类

本章所讨论的文本分类方法主要是对文档对象进行分类。本章从文本分类的基本概念、应用场景及分类特征优化等方面对文本分类的技术进行了系统的介绍。本章介绍了三类重要的分类模型：朴素贝叶斯模型、向量空间模型、支持向量机模型。

## 第 5 章 文本聚类

本章介绍了对文档对象进行聚类描述的主要技术方法，主要涵盖了扁平式聚类和凝聚式聚类两大基本问题解决思路。本章还介绍了如何对聚类结果进行分析，以及对聚类的特征进行优化等相关内容。对特殊文本对象的聚类技术的介绍也是本章的重点内容，具体包括半监督聚类、短文本聚类及流数据聚类。

## 第 6 章 序列标注

序列标注是特殊的分类问题，很多文本分析任务都需要抽象成序列标注问题进行解决。本章介绍了当前三类重要的序列标注基础模型，即隐马尔可夫模型、最大熵马尔可夫模型及条件随机场。本章还介绍了各模型的主要特征、优点和缺点，并提供了具体的应用范例。

## 第 7 章 信息检索

本章介绍了如何根据用户的特定信息需求，从在线环境中有效地提取重要的文本对象并进行反馈。除了介绍信息检索的重要应用场景，本章还讨论了三类主流的模型方案：基于空间模型的信息检索、基于概率模型的信息检索、基于语言模型的信息检索。

## 第 8 章 文本摘要

本章介绍了如何基于已有文本内容对信息进行压缩，并从中提取有价值的、关键的文本要素。文本摘要技术包括关键词提取和关键句提取，前者是本章介绍的重点。本章还介绍了很多经典的对词汇的关键词进行量化评估的指标，同时介绍了当前主流的基于图模型的关键词提取算法。

## 第 9 章 口碑分析

本章介绍了如何从在线平台的用户评论文本数据中提取有价值的产品信息。一方面，本章讨论了如何通过词典或语料集合对在线评价对象进行提取；另一方面，本章介绍了如何在不同的粒度水平上挖掘用户对于产品或服务的情感态度。

## 第 10 章 社交网络分析

社交网络是重要的互联网应用场景。本章介绍了很多社交网络上的文本分析

任务及具体的技术方案，包括社交网络的虚拟社区发现、用户影响力分析、情感分析、话题发现与演化，以及信息检索。本章还介绍了如何将社交网络的多属性特征和图结构特征有机地结合到文本分析技术框架中。

## 第 11 章 深度学习与 NLP

本章介绍了当前热门的深度学习技术在文本分析中的应用。深度学习以神经网络为基础模型。本章分别介绍了基于多层感知器模型和循环神经网络的深度学习文本分析技术。对于循环神经网络，本章特别介绍了词嵌入模型和机器翻译技术。

## 第 12 章 实证研究

本章介绍了文本分析技术在互联网领域中的管理类应用，讲述了如何通过实证研究来挖掘在线平台上的用户行为，并结合研究结果有针对性地提供管理决策建议。本章还介绍了文本分析技术在互联网医疗中的具体应用，以真实的场景、数据为依托，为从事互联网运营相关工作的读者提供了有价值的解决问题的思路。

## 第 13 章 总结

作为结束语，本章简要回顾了全书的核心内容，并为文本分析领域的工作者提供了若干条有价值的实践经验。

## 本书读者对象

- 从事数据分析、文本分析相关职业的技术人员、网络运营人员；
- 所学专业与计算机技术、互联网技术、语言学相关的本科生及研究生；
- 计算机科学、自然语言处理等领域的大学教师及科研工作者；
- 其他对文本分析有兴趣爱好的人员。

# 目 录

<b>第1章 概述 .....</b>	1
1.1 网络运营与文本分析.....	1
1.1.1 互联网运营的战略思维 .....	1
1.1.2 网络运营与大数据文本分析 .....	2
1.2 文本分析的4V特征.....	4
1.2.1 Volume 特征.....	4
1.2.2 Variety 特征 .....	5
1.2.3 Value 特征 .....	6
1.2.4 Velocity 特征.....	7
1.3 在线文本分析应用 .....	8
1.3.1 在线文本分析的管理类应用 .....	9
1.3.2 在线文本分析的内容类应用 .....	12
1.4 本章小结 .....	16
<b>第2章 预备知识.....</b>	18
2.1 文本挖掘的主要任务.....	18
2.2 语义分析与语法分析.....	20
2.3 文本的结构化分析 .....	21
2.4 文本的标准化分析 .....	24
2.5 机器学习的基本概念.....	24
2.5.1 机器学习与深度学习 .....	25
2.5.2 机器学习的基本要素 .....	33
2.6 机器学习的重要问题.....	36
2.6.1 概率图模型 .....	36

2.6.2 判别式模型和产生式模型 .....	39
2.6.3 机器学习模型求解 .....	40
2.6.4 模型过拟合 .....	43
2.7 本章小结 .....	45

### 第3章 文本建模..... 46

3.1 文本建模的基本概念 .....	46
3.2 文本建模的应用场景 .....	48
3.2.1 主体角色识别 .....	48
3.2.2 语言风格分析 .....	49
3.2.3 智能系统 .....	49
3.2.4 文本表示 .....	50
3.2.5 文本降维 .....	50
3.2.6 话题分析 .....	50
3.3 语言学建模概述 .....	51
3.4 词标注分析 .....	52
3.5 句法分析 .....	55
3.5.1 转换生成语法 .....	56
3.5.2 依存句法 .....	56
3.6 知识库与语义网 .....	58
3.7 统计学建模概述 .....	59
3.8 向量空间模型 .....	61
3.9 LSI 模型 .....	64
3.9.1 SVD .....	64
3.9.2 基于 SVD 的降维分析 .....	66
3.10 Unigram 模型 .....	67
3.11 pLSI 模型 .....	67
3.11.1 pLSI 的模型结构 .....	67
3.11.2 pLSI 的参数估计 .....	68

3.12 LDA 主题模型 .....	70
3.12.1 LDA 的模型结构 .....	70
3.12.2 LDA 的参数估计 .....	72
3.13 主题模型拓展 .....	75
3.13.1 相关主题模型 .....	76
3.13.2 层次主题模型 .....	77
3.13.3 动态主题模型 .....	80
3.13.4 句子主题模型 .....	82
3.14 基于词汇的统计学建模方法 .....	83
3.15 本章小结 .....	86

## 第4章 文本分类 ..... 88

4.1 文本分类的基本概念 .....	88
4.2 文本分类的应用场景 .....	89
4.2.1 文档有用性判断 .....	89
4.2.2 口碑情感分析 .....	90
4.2.3 负面信息识别 .....	90
4.2.4 信息检索 .....	90
4.3 朴素贝叶斯模型 .....	91
4.3.1 贝努利模型 .....	91
4.3.2 多项式模型 .....	93
4.3.3 模型参数平滑 .....	94
4.4 向量空间模型 .....	95
4.4.1 Rocchio 方法 .....	95
4.4.2 KNN 方法 .....	96
4.5 SVM 模型 .....	97
4.5.1 硬间隔 SVM .....	97
4.5.2 软间隔 SVM .....	100
4.6 文本分类的评价 .....	102

4.6.1 二元分类评价 .....	102
4.6.2 多类问题评价 .....	104
4.6.3 分类测试集 .....	105
4.7 分类特征优化 .....	106
4.7.1 分类特征提取 .....	106
4.7.2 分类特征转化 .....	112
4.7.3 分类特征扩展 .....	114
4.8 分类学习策略优化 .....	117
4.8.1 AdaBoost 算法 .....	117
4.8.2 主动式学习 .....	118
4.8.3 迁移学习 .....	119
4.9 本章小结 .....	119

## 第5章 文本聚类 ..... 121

5.1 文本聚类的基本概念 .....	121
5.2 文本聚类的应用场景 .....	122
5.2.1 探索分析 .....	122
5.2.2 降维 .....	123
5.2.3 信息检索 .....	123
5.3 扁平式聚类 .....	124
5.3.1 K-均值算法 .....	125
5.3.2 基于模型的聚类 .....	128
5.4 凝聚式聚类 .....	132
5.4.1 层次聚类 .....	132
5.4.2 基于簇距离的聚类过程 .....	132
5.4.3 算法停止条件 .....	135
5.5 聚类结果分析 .....	136
5.5.1 聚类算法评估 .....	136
5.5.2 聚类标签生成 .....	138

5.6 聚类特征优化 .....	140
5.6.1 基于迭代的方法 .....	141
5.6.2 无监督指标 .....	141
5.7 半监督聚类 .....	143
5.7.1 迁移学习 .....	144
5.7.2 AP 算法 .....	145
5.8 短文本聚类 .....	146
5.8.1 文本特征补充 .....	146
5.8.2 TermCut 算法 .....	148
5.8.3 Dirichlet 多项式混合模型 .....	149
5.9 流数据聚类 .....	151
5.9.1 OSKM 算法 .....	151
5.9.2 可拓展 K-means 算法 .....	152
5.10 本章小结 .....	153

<b>第6章 序列标注.....</b>	<b>155</b>
6.1 序列标注的基本概念 .....	155
6.2 序列标注的应用场景 .....	157
6.2.1 词性标注 .....	157
6.2.2 命名实体识别 .....	157
6.2.3 分词 .....	157
6.3 HMM .....	158
6.3.1 HMM 的概率计算问题 .....	160
6.3.2 HMM 的学习问题 .....	162
6.3.3 HMM 的预测问题 .....	164
6.4 最大熵模型和最大熵马尔可夫模型 .....	166
6.4.1 最大熵模型 .....	167
6.4.2 最大熵马尔可夫模型 .....	170
6.5 条件随机场 .....	172

6.5.1 标注偏置问题.....	172
6.5.2 条件随机场的基本原理.....	174
6.6 本章小结 .....	176

## 第7章 信息检索..... 177

7.1 信息检索的基本概念.....	177
7.2 信息检索的应用场景.....	180
7.2.1 搜索引擎 .....	180
7.2.2 内容推荐 .....	182
7.3 基于空间模型的信息检索.....	184
7.3.1 文档查找 .....	184
7.3.2 文档排序 .....	185
7.3.3 系统评价 .....	187
7.4 基于概率模型的信息检索.....	190
7.4.1 二值独立模型 .....	191
7.4.2 模型参数估计 .....	193
7.5 基于语言模型的信息检索.....	196
7.5.1 语言模型 .....	196
7.5.2 查询似然模型 .....	198
7.6 本章小结 .....	201

## 第8章 文本摘要..... 203

8.1 文本摘要的基本概念.....	203
8.2 文本摘要的应用场景.....	206
8.2.1 信息检索 .....	206
8.2.2 信息压缩 .....	207
8.2.3 用户画像 .....	208
8.2.4 知识管理 .....	209
8.3 关键词提取的特征设计.....	210

8.3.1 词频特征 .....	210
8.3.2 词汇基础特征 .....	211
8.3.3 词汇位置特征 .....	212
8.3.4 词汇标记特征 .....	214
8.4 关键词提取的有监督算法 .....	214
8.5 关键词提取的无监督算法 .....	217
8.5.1 简单指标设计 .....	217
8.5.2 复合指标设计 .....	217
8.6 基于图模型的关键词提取算法 .....	218
8.6.1 图模型静态指标算法 .....	220
8.6.2 图模型动态指标算法 .....	223
8.7 关键词提取的技术优化 .....	226
8.7.1 长文本问题优化 .....	227
8.7.2 短文本问题优化 .....	228
8.7.3 多主题特征优化 .....	229
8.7.4 时序特征优化 .....	232
8.7.5 歧义问题优化 .....	233
8.8 关键短语提取 .....	234
8.8.1 短语性指标 .....	235
8.8.2 信息性指标 .....	235
8.9 关键句提取 .....	236
8.9.1 基于词汇关键性的方法 .....	236
8.9.2 基于句子特征的方法 .....	237
8.9.3 基于图模型的方法 .....	238
8.10 本章小结 .....	240
<b>第9章 口碑分析 .....</b>	<b>241</b>
9.1 口碑分析的基本概念 .....	242
9.2 口碑分析的应用场景 .....	243