

Text Data Mining

Chengqing Zong • Rui Xia • Jiajun Zhang

Text Data Mining

Chengqing Zong
Institute of Automation
Chinese Academy of Sciences
Beijing, Beijing, China

Rui Xia
School of Computer Science & Engineering
Nanjing University of Science
and Technology
Nanjing, Jiangsu, China

Jiajun Zhang
Institute of Automation
Chinese Academy of Sciences
Beijing, Beijing, China

ISBN 978-981-16-0099-9 ISBN 978-981-16-0100-2 (eBook)
<https://doi.org/10.1007/978-981-16-0100-2>

Jointly published with Tsinghua University Press

The print edition is not for sale in China (Mainland). Customers from China (Mainland) please order the print book from: Tsinghua University Press.

© Tsinghua University Press 2021

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed.

The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

The publisher, the authors, and the editors are safe to assume that the advice and information in this book are believed to be true and accurate at the date of publication. Neither the publisher nor the authors or the editors give a warranty, expressed or implied, with respect to the material contained herein or for any errors or omissions that may have been made. The publisher remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

This Springer imprint is published by the registered company Springer Nature Singapore Pte Ltd.
The registered company address is: 152 Beach Road, #21-01/04 Gateway East, Singapore 189721, Singapore

Foreword

We are living in the Big Data era. Over 80% of real-world data are unstructured in the form of natural language text, such as books, news reports, research articles, social media messages, and webpages. Although data mining and machine learning have been popular in data analysis, most data mining methods handle only structured or semistructured data. In comparison with mining structured data, mining unstructured text data is more challenging but will also play a more essential role in turning massive data into structured knowledge. It is no wonder we have witnessed such a dramatic upsurge in the research on text mining and natural language processing and their applications in recent years.

Text mining is a confluence of natural language processing, data mining, machine learning, and statistics used to mine knowledge from unstructured text. There have already been multiple textbooks dedicated to data mining, machine learning, statistics, and natural language processing. However, we seriously lack textbooks on text mining that systematically introduce important topics and up-to-date methods for text mining. This book, “Text Data Mining,” bridges this gap nicely. It is the first textbook, and a brilliant one, on text data mining that not only introduces foundational issues but also offers comprehensive and state-of-the-art coverage of the important and ongoing research themes on text mining. The in-depth treatment of a wide spectrum of text mining themes and clear introduction to the state-of-the-art deep learning methods for text mining make the book unique, timely, and authoritative. It is a great textbook for graduate students as well as a valuable handbook for practitioners working on text mining, natural language processing, data mining, and machine learning and their applications. This book is written by three pioneering researchers and highly reputed experts in the fields of natural language processing and text mining. The first author has written an authoritative and popular textbook on natural language processing that has been adopted as a standard textbook for university undergraduate and first-year graduate students in China. However, this new text mining book has completely different coverage from his NLP textbook and offers new and complementary text mining themes. Both books can be studied independently, although I would strongly encourage students working on NLP and text mining to study both.

This text mining book starts with text preprocessing, including both English and Chinese text preprocessing, and proceeds to text representation, covering the vector space model and distributed representation of words, phrases, sentences, and documents, in both statistical modeling and deep learning models. It then introduces feature selection methods, statistical learning methods, and deep neural network methods, including multilayer feed-forward neural networks, convolutional neural networks, and recurrent neural networks, for document classification. It next proceeds to text clustering, covering sample and cluster similarities, various clustering methods, and clustering evaluation. After introducing the fundamental theories and methods of text mining, the book uses five chapters to cover a wide spectrum of text mining applications, including topic modeling (which is also treated as a fundamental issue from some viewpoints but can be used independently), sentiment analysis and opinion mining, topic detection and tracking, information extraction, and automated document summarization. These themes are active research frontiers in text mining and are covered comprehensively and thoroughly, with a good balance between classical methods and recent developments, including deep learning methods.

As a data mining researcher, I have recently been deeply involved in text mining due to the need to handle the large scale of real-world data. I could not find a good text mining textbook written in English or Chinese to learn and teach. It is exciting to see that this book provides such a comprehensive and cutting-edge introduction. I believe this book will benefit data science researchers, graduate students, and those who want to include text mining in practical applications. I loved reading this book and recommend it highly to everyone who wants to learn text mining!

ACM Fellow and IEEE Fellow
Michael Aiken Chair Professor
Department of Computer Science
University of Illinois at Urbana-Champaign
Champaign, IL, USA

Jiawei Han



Preface

With the rapid development and popularization of Internet and mobile communication technologies, text data mining has attracted much attention. In particular, with the wide use of new technologies such as cloud computing, big data, and deep learning, text mining has begun playing an increasingly important role in many application fields, such as opinion mining and medical and financial data analysis, showing broad application prospects.

Although I was supervising graduate students studying text classification and automatic summarization more than ten years ago, I did not have a clear understanding of the overall concept of text data mining and only regarded the research topics as specific applications of natural language processing. Professor Jiawei Han's book *Data Mining: Concepts and Technology*, published by Elsevier, Professor Bing Liu's *Web Data Mining*, published by Springer, and other books have greatly benefited me. Every time I listen to their talks and discuss these topics with them face to face, I have benefited immensely. I was inspired to write this book for the course Text Data Mining, which I was invited to teach to graduates of the University of Chinese Academy of Sciences. At the end of 2015, I accepted the invitation and began to prepare the content design and selection of materials for the course. I had to study a large number of related papers, books, and other materials and began to seriously think of the rich connotation and extension of the term Text Data Mining. After more than a year's study, I started to compile the courseware. With teaching practice, the outline of the concept has gradually formed.

Rui Xia and Jiajun Zhang, two talented young people, helped me materialize my original writing plan. Rui Xia received his master's degree in 2007 and was admitted to the Institute of Automation, Chinese Academy of Sciences, and studied for Ph.D. degree under my supervision. He was engaged in sentiment classification and took it as the research topic of his Ph.D. dissertation. After he received his Ph.D. degree in 2011, his interests extended to opinion mining, text clustering and classification, topic modeling, event detection and tracking, and other related topics. He has published a series of influential papers in the field of sentiment analysis and opinion mining. He received the ACL 2019 outstanding paper award, and his paper on ensemble learning for sentiment classification has been cited more than

600 times. Jiajun Zhang joined our institute after he graduated from university in 2006 and studied in my group in pursuit of his Ph.D. degree. He mainly engaged in machine translation research, but he performed well in many research topics, such as multilanguage automatic summarization, information extraction, and human–computer dialogue systems. Since 2016, he has been teaching some parts of the course on Natural Language Processing in cooperation with me, such as machine translation, automatic summarization, and text classification, at the University of Chinese Academy of Sciences; this course is very popular with students. With the solid theoretical foundation of these two talents and their keen scientific insights, I am gratified that many cutting-edge technical methods and research results could be verified and practiced and included in this book.

From early 2016 to June 2019, when the Chinese version of this book was published, it took more than three years. In these three years, most holidays, weekends, and other spare times of ours were devoted to the writing of this book. It was really suffering to endure the numerous modifications or even rewriting, but we were also very happy. We started to translate the Chinese version into English in the second half of 2019. Some more recent topics, including BERT (bidirectional encoder representations from transformers), have been added to the English version.

As a cross domain of natural language processing and machine learning, text data mining faces the double challenges of the two domains and has broad application to the Internet and equipment for mobile communication. The topics and techniques presented in this book are all the technical foundations needed to develop such practical systems and have attracted much attention in recent years. It is hoped that this book will provide a comprehensive understanding for students, professors, and researchers in related areas. However, I must admit that due to the limitation of the authors' ability and breadth of knowledge, as well as the lack of time and energy, there must be some omissions or mistakes in the book. We will be very grateful if readers provide criticism, corrections, and any suggestions.

Beijing, China
20 May 2020

Chengqing Zong

Acknowledgments

During the writing process and after the completion of the first draft of the Chinese version of this book, many experts from related fields reviewed selected chapters and gave us valuable comments and suggestions. They are (in alphabetical order) Xianpei Han, Yu Hong, Shoushan Li, Kang Liu, Xiaojun Wan, Kang Xu, Chengzhi Zhang, and Xin Zhao. In addition, we also received help from several graduate students and Ph.D. candidates (also in alphabetical order): Hongjie Cai, Zixiang Ding, Huihui He, Xiao Jin, Junjie Li, Mei Li, Yuchen Liu, Cong Ma, Liqun Ma, Xiangqing Shen, Jingyuan Sun, Fanfan Wang, Leyi Wang, Qain Wang, Weikang Wang, Yining Wang, Kaizhou Xuan, Shiliang Zheng, and Long Zhou. They helped us to double check and confirm English expressions, references, and web addresses and to redraw several charts in the book, which saved us much time. We would like to express our heartfelt thanks to all of them!

We also want to sincerely thank Professor Jiawei Han for his guidance and suggestions on this book. We are honored that he was willing to write the foreword to this book despite his busy schedule. Finally, we would like to recognize Ms. Hui Xue and Qian Wang, Tsinghua University Press, and Ms. Celine Chang, and Ms. Suraj Kumar, Springer, for their great help!

Beijing, China
Nanjing, China
Beijing, China
20 May 2020

Chengqing Zong
Rui Xia
Jiajun Zhang

Contents

- 1 Introduction** 1
 - 1.1 The Basic Concepts 1
 - 1.2 Main Tasks of Text Data Mining 3
 - 1.3 Existing Challenges in Text Data Mining 6
 - 1.4 Overview and Organization of This Book 9
 - 1.5 Further Reading 12
- 2 Data Annotation and Preprocessing** 15
 - 2.1 Data Acquisition..... 15
 - 2.2 Data Preprocessing 20
 - 2.3 Data Annotation 22
 - 2.4 Basic Tools of NLP 25
 - 2.4.1 Tokenization and POS Tagging 25
 - 2.4.2 Syntactic Parser 27
 - 2.4.3 *N*-gram Language Model 29
 - 2.5 Further Reading 30
- 3 Text Representation** 33
 - 3.1 Vector Space Model 33
 - 3.1.1 Basic Concepts..... 33
 - 3.1.2 Vector Space Construction 34
 - 3.1.3 Text Length Normalization..... 36
 - 3.1.4 Feature Engineering 37
 - 3.1.5 Other Text Representation Methods 39
 - 3.2 Distributed Representation of Words 40
 - 3.2.1 Neural Network Language Model 41
 - 3.2.2 C&W Model 45
 - 3.2.3 CBOW and Skip-Gram Model..... 47
 - 3.2.4 Noise Contrastive Estimation and Negative Sampling.... 49
 - 3.2.5 Distributed Representation Based on the Hybrid
Character-Word Method..... 51

3.3	Distributed Representation of Phrases	53
3.3.1	Distributed Representation Based on the Bag-of-Words Model	54
3.3.2	Distributed Representation Based on Autoencoder	54
3.4	Distributed Representation of Sentences	58
3.4.1	General Sentence Representation	59
3.4.2	Task-Oriented Sentence Representation	63
3.5	Distributed Representation of Documents	66
3.5.1	General Distributed Representation of Documents	67
3.5.2	Task-Oriented Distributed Representation of Documents	69
3.6	Further Reading	72
4	Text Representation with Pretraining and Fine-Tuning	75
4.1	ELMo: Embeddings from Language Models	75
4.1.1	Pretraining Bidirectional LSTM Language Models.....	76
4.1.2	Contextualized ELMo Embeddings for Downstream Tasks.....	77
4.2	GPT: Generative Pretraining.....	78
4.2.1	Transformer	78
4.2.2	Pretraining the Transformer Decoder	80
4.2.3	Fine-Tuning the Transformer Decoder	81
4.3	BERT: Bidirectional Encoder Representations from Transformer.....	82
4.3.1	BERT: Pretraining	83
4.3.2	BERT: Fine-Tuning.....	86
4.3.3	XLNet: Generalized Autoregressive Pretraining	86
4.3.4	UniLM	89
4.4	Further Reading	90
5	Text Classification	93
5.1	The Traditional Framework of Text Classification	93
5.2	Feature Selection	95
5.2.1	Mutual Information	96
5.2.2	Information Gain	99
5.2.3	The Chi-Squared Test Method	100
5.2.4	Other Methods	101
5.3	Traditional Machine Learning Algorithms for Text Classification	102
5.3.1	Naïve Bayes.....	103
5.3.2	Logistic/Softmax and Maximum Entropy	105
5.3.3	Support Vector Machine.....	107
5.3.4	Ensemble Methods	110

5.4	Deep Learning Methods	111
5.4.1	Multilayer Feed-Forward Neural Network	111
5.4.2	Convolutional Neural Network	113
5.4.3	Recurrent Neural Network	115
5.5	Evaluation of Text Classification	120
5.6	Further Reading	123
6	Text Clustering	125
6.1	Text Similarity Measures	125
6.1.1	The Similarity Between Documents	125
6.1.2	The Similarity Between Clusters	128
6.2	Text Clustering Algorithms	129
6.2.1	K -Means Clustering	129
6.2.2	Single-Pass Clustering	133
6.2.3	Hierarchical Clustering	136
6.2.4	Density-Based Clustering	138
6.3	Evaluation of Clustering	141
6.3.1	External Criteria	141
6.3.2	Internal Criteria	142
6.4	Further Reading	143
7	Topic Model	145
7.1	The History of Topic Modeling	145
7.2	Latent Semantic Analysis	146
7.2.1	Singular Value Decomposition of the Term-by-Document Matrix	147
7.2.2	Conceptual Representation and Similarity Computation	148
7.3	Probabilistic Latent Semantic Analysis	150
7.3.1	Model Hypothesis	150
7.3.2	Parameter Learning	151
7.4	Latent Dirichlet Allocation	153
7.4.1	Model Hypothesis	153
7.4.2	Joint Probability	155
7.4.3	Inference in LDA	158
7.4.4	Inference for New Documents	160
7.5	Further Reading	161
8	Sentiment Analysis and Opinion Mining	163
8.1	History of Sentiment Analysis and Opinion Mining	163
8.2	Categorization of Sentiment Analysis Tasks	164
8.2.1	Categorization According to Task Output	164
8.2.2	According to Analysis Granularity	165

8.3	Methods for Document/Sentence-Level Sentiment Analysis	168
8.3.1	Lexicon- and Rule-Based Methods	169
8.3.2	Traditional Machine Learning Methods	170
8.3.3	Deep Learning Methods	174
8.4	Word-Level Sentiment Analysis and Sentiment Lexicon Construction	178
8.4.1	Knowledgebase-Based Methods	178
8.4.2	Corpus-Based Methods	179
8.4.3	Evaluation of Sentiment Lexicons	182
8.5	Aspect-Level Sentiment Analysis	183
8.5.1	Aspect Term Extraction	183
8.5.2	Aspect-Level Sentiment Classification	186
8.5.3	Generative Modeling of Topics and Sentiments	191
8.6	Special Issues in Sentiment Analysis	193
8.6.1	Sentiment Polarity Shift	193
8.6.2	Domain Adaptation	195
8.7	Further Reading	198
9	Topic Detection and Tracking	201
9.1	History of Topic Detection and Tracking	201
9.2	Terminology and Task Definition	202
9.2.1	Terminology	202
9.2.2	Task	203
9.3	Story/Topic Representation and Similarity Computation	206
9.4	Topic Detection	209
9.4.1	Online Topic Detection	209
9.4.2	Retrospective Topic Detection	211
9.5	Topic Tracking	212
9.6	Evaluation	213
9.7	Social Media Topic Detection and Tracking	215
9.7.1	Social Media Topic Detection	216
9.7.2	Social Media Topic Tracking	217
9.8	Bursty Topic Detection	217
9.8.1	Burst State Detection	218
9.8.2	Document-Pivot Methods	221
9.8.3	Feature-Pivot Methods	222
9.9	Further Reading	224
10	Information Extraction	227
10.1	Concepts and History	227
10.2	Named Entity Recognition	229
10.2.1	Rule-based Named Entity Recognition	230
10.2.2	Supervised Named Entity Recognition Method	231
10.2.3	Semisupervised Named Entity Recognition Method	239
10.2.4	Evaluation of Named Entity Recognition Methods	241

10.3	Entity Disambiguation	242
10.3.1	Clustering-Based Entity Disambiguation Method	243
10.3.2	Linking-Based Entity Disambiguation	248
10.3.3	Evaluation of Entity Disambiguation	254
10.4	Relation Extraction	256
10.4.1	Relation Classification Using Discrete Features	258
10.4.2	Relation Classification Using Distributed Features	265
10.4.3	Relation Classification Based on Distant Supervision	268
10.4.4	Evaluation of Relation Classification	269
10.5	Event Extraction	270
10.5.1	Event Description Template	270
10.5.2	Event Extraction Method	272
10.5.3	Evaluation of Event Extraction	281
10.6	Further Reading	281
11	Automatic Text Summarization	285
11.1	Main Tasks in Text Summarization	285
11.2	Extraction-Based Summarization	287
11.2.1	Sentence Importance Estimation	287
11.2.2	Constraint-Based Summarization Algorithms	298
11.3	Compression-Based Automatic Summarization	299
11.3.1	Sentence Compression Method	300
11.3.2	Automatic Summarization Based on Sentence Compression	305
11.4	Abstractive Automatic Summarization	307
11.4.1	Abstractive Summarization Based on Information Fusion	307
11.4.2	Abstractive Summarization Based on the Encoder-Decoder Framework	313
11.5	Query-Based Automatic Summarization	316
11.5.1	Relevance Calculation Based on the Language Model ...	317
11.5.2	Relevance Calculation Based on Keyword Co-occurrence	317
11.5.3	Graph-Based Relevance Calculation Method	318
11.6	Crosslingual and Multilingual Automatic Summarization	319
11.6.1	Crosslingual Automatic Summarization	319
11.6.2	Multilingual Automatic Summarization	323
11.7	Summary Quality Evaluation and Evaluation Workshops	325
11.7.1	Summary Quality Evaluation Methods	325
11.7.2	Evaluation Workshops	330
11.8	Further Reading	332
	References	335

About the Authors

Chengqing Zong is a Professor at the National Laboratory of Pattern Recognition (NLPR), Institute of Automation, Chinese Academy of Sciences (CASIA) and an adjunct professor in the School of Artificial Intelligence at University of Chinese Academy of Sciences (UCAS). He authored the book “Statistical Natural Language Processing” (which is in Chinese and sold more than 32K copies) and has published more than 200 papers on machine translation, natural language processing, and cognitive linguistics. He served as the chair for numerous prestigious conferences, such as ACL, COLING, AACL, and IJCAI, and has served as an associate editor for journals, such as ACM TALLIP and ACTA Automatic Sinica, and as an editorial board member for journals, including IEEE Intelligent Systems, Journal of Computer Science and Technology, and Machine Translation. He is currently the President of the Asian Federation of Natural Language Processing (AFNLP) and a member of the International Committee on Computational Linguistics (ICCL).

Rui Xia is a Professor at the School of Computer Science and Engineering, Nanjing University of Science and Technology, China. He has published more than 50 papers in high-quality journals and top-tiered conferences in the field of natural language processing and text data mining. He serves as area chair and senior program committee member for several top conferences, such as EMNLP, COLING, IJCAI, and AACL. He received the outstanding paper award of ACL 2019, and the Distinguished Young Scholar award from the Natural Science Foundation of Jiangsu Province, China in 2020.

Jiajun Zhang is a Professor at the NLPR, CASIA and an adjunct professor in the SAIU of UCAS. He has published more than 80 conference papers and journal articles on natural language processing and text mining and received 5 best paper awards. He served as the area chair or on the senior program committees for several top conferences, such as ACL, EMNLP, COLING, AACL, and IJCAI. He is the

deputy director of China's Machine Translation Technical Committee of the Chinese Information Processing Society of China. He received the Qian Wei-Chang Science and Technology Award of Chinese Information Processing and the CIPS Hanvon Youth Innovation Award. He was supported by the Elite Scientists Sponsorship Program of China Association for Science and Technology (CAST).

Acronyms

ACE	Automatic content extraction
AMR	Abstract meaning representation
ATT	Adaptive topic tracking
AUC	Area under the ROC curve
Bagging	Bootstrap aggregating
BERT	Bidirectional encoder representations from transformer
BFGS	Broyden–Fletcher–Goldfarb–Shanno
Bi-LSTM	Bidirectional long short-term memory
BIO	Begin–inside–outside
BLEU	Bilingual evaluation understudy
BOW	Bag of words
BP	Back-propagation
BPTS	Back-propagation through structure
BPTT	Back-propagation through time
BRAE	Bilingually constrained recursive autoencoder
CBOW	Continuous bag-of-words
CFG	Context-free grammar
CNN	Convolutional neural network
CRF	Conditional random field
C&W	Collobert and Weston
CWS	Chinese word segmentation
DBSCAN	Density-based spatial clustering of applications with noise
DF	Document frequency
DL	Deep learning
DMN	Deep memory network
ELMo	Embeddings from language models
EM	Expectation maximization
EP	Expectation propagation
FAR	False alarm rate
FNN	Feed-forward neural network
GPT	Generative pretraining

GRU	Gated recurrent unit
HAC	Hierarchical agglomerative clustering
HMM	Hidden Markov model
HTD	Hierarchical topic detection
ICA	Independent component analysis
IDF	Inverse document frequency
IE	Information extraction
IG	Information gain
KBP	Knowledge base population
KDD	Knowledge discovery in databases
KKT	Karush–Kuhn–Tucker
K-L	Kullback–Leibler
L-BFGS	Limited-memory Broyden–Fletcher–Goldfarb–Shanno
LDA	Latent Dirichlet allocation
LM	Language model
LSA	Latent semantic analysis
LSI	Latent semantic indexing
LSTM	Long short-term memory
MCMC	Markov chain Monte Carlo
MDR	Missed detection rate
ME	Maximum entropy
MI	Mutual information
ML	Machine learning
MLE	Maximum likelihood estimation
MMR	Maximum marginal relevance
MUC	Message understanding conference
NB	Naïve Bayes
NCE	Noise contrastive estimation
NED	New event detection
NER	Named entity recognition
NLP	Natural language processing
NNLM	Neural network language model
PCA	Principal component analysis
PLSA	Probabilistic latent semantic analysis
PLSI	Probabilistic latent semantic indexing
PMI	Pointwise mutual information
P-R	Precision–recall
PU	Positive-unlabeled
PCFG	Probabilistic context-free grammar
PMI-IR	Pointwise mutual information—information retrieval
POS	Part of speech
PV-DBOW	Distributed bag-of-words version of the paragraph vector
PV-DM	Paragraph vector with sentence as distributed memory
Q&A	Question and answering
RAE	Recursive autoencoder

RecurNN	Recursive neural network
RG	Referent graph
RNN	Recurrent neural network
ROC	Receiver operating characteristic
ROUGE	Recall-oriented understudy for gisting evaluation
RTD	Retrospective topic detection
SCL	Structure correspondence learning
SCU	Summary content unit
SMO	Sequential minimal optimization
SO	Semantic orientation
SRL	Semantic role labeling
SS	Story segmentation
SST	Stanford sentiment treebank
STC	Suffix tree clustering
SVD	Singular value decomposition
SVM	Support vector machine
TAC	Text analysis conference
TD	Topic detection
TDT	Topic detection and tracking
TF	Term frequency
TF-IDF	Term frequency—inverted document frequency
UGC	User-generated context
UniLM	Unified pretraining language model
VBEM	Variational Bayes expectation maximization
VSM	Vector space model
WCSS	Within-cluster sum of squares
WSD	Word sense disambiguation