

《自然语言处理：大模型理论与实践》

赵宇 任福继 陈星延 陈中普 陈珍珠 施龙 张阳 李庆 谢志龙

2024 年

序 言

随着人工智能技术的飞速发展，自然语言处理成为了计算机科学与人工智能领域中不可或缺的关键技术之一。作为一名长期致力于人工智能和自然语言处理研究的学者，我见证了这一领域的迅猛变革，从基于规则的方法到现今基于深度学习的革命性技术，尤其是大模型技术的应用，给自然语言处理领域带来了前所未有的机遇和挑战。

记得 1988-1990 年，我在博士期间开发了一个日中机器翻译系统，用的是 Fortran 语言，在汉字输入输出、字符串处理上费尽了脑筋。去年，我有机会应邀在很多场合演讲大模型和人工智能，提到过目前的大模型在理论上并无多大创新。我列举了国际研究同行发来的我们 1999 年发表的论文《Automatic Derivation of Programs for Image Processing from Natural Language Descriptions》（从自然语言描述自动推导程序），尽管由于周知的数据规模和算力限制，但模型系统成功的理论基础与大模型极其一致。既然目前大模型理论上并无多大创新，那为何现今它具有如此巨大的能力，推动人工智能奔腾向前？我认为就是本书要提到的涌现。

本书以其独特的视角和结构，全面系统地介绍了大模型技术在自然语言处理中的理论与实践。全书内容丰富，涵盖了语言模型的基础知识、大模型的关键技术以及大模型在实际中的应用实践，不仅为初学者提供了详尽的入门指南，也为研究人员和专业人士提供了深入的技术解析和实用的开发案例。

本书从自然语言处理的背景知识出发，逐步引入词向量、统计语言模型、神经语言模型和预训练语言模型，然后详细介绍了大模型的架构、训练方法、应用及评估策略，并通过丰富的实例和习题，帮助读者加深理解和掌握大模型技术的核心内容。特别的是，本书对大模型的本地开发和应用开发部分，提供了实战演练和实用的代码示例，对实际开发工作具有重要的指导意义。

本书的出版正值大模型技术蓬勃发展的关键时期，对于高校本科生、研究生、教学科研人员，以及从事自然语言处理研究和开发的专业人士来说，都是一本不可多得的参考书籍。通过阅读本书，读者不仅可以系统地掌握自然语言

处理的基础知识和前沿技术，还能在实际操作中提升自身的开发和研究能力，探索自然语言处理的广阔前景。

我深信，随着自然语言处理技术的不断发展和应用，未来将会有更多的突破和创新，推动人类与机器之间更加自然和智能的交流方式，开创人工智能更加辉煌的未来。

任福继

日本工程院院士、中国人工智能学会理事会名誉理事长

2024 年 6 月于成都

前言

自然语言处理 (Natural Language Processing, NLP) 是计算机科学与人工智能交叉领域中的一门关键技术, 其目标是使计算机能够理解、解释、生成人类语言。在当今人工智能时代, NLP 技术已经深刻地渗透到我们日常生活的方方面面, 从智能助手、语音识别到机器翻译和文本生成, NLP 正以惊人的速度改变着我们的生活方式。特别的是, 2022 年底以 ChatGPT 为代表的大模型技术横空出世, 进一步推动了新一代人工智能技术的发展。大模型技术颠覆了自然语言处理领域传统的知识体系。然而, 目前以大模型技术为主线介绍自然语言处理知识的教材较为缺乏。基于此考虑, 催生了我编著本教材的想法。

本教材以自然语言处理中语言模型为主线, 主要内容分为三部分, 包括语言模型基础、大语言模型关键技术和大模型实践。在此之前, 首先介绍了自然语言处理的背景知识。然后, 在语言模型基础部分介绍了词向量、统计语言模型、神经语言模型和预训练语言模型。接着, 在大模型理论部分介绍大模型的架构、训练、利用与评估等。最后, 在大模型实践部分介绍了大模型的本地开发和应用开发等。

本教材主要针对高校本科生、研究生以及教学科研人员, 作为教学用书。当然, 也适用于计算语言学家、语言学家、数据科学家和 NLP 开发人员等专业人士。考虑到不同读者的学科差异, 本书在附录部分介绍了概率论、信息论、机器学习与强化学习等 NLP 交叉学科的基础知识。阅读本教材最好具备 Python 的编程知识。

在写作本教材的过程中, 我深切地感受到自然语言处理的迅猛发展。从传统的基于规则的方法到现今基于深度学习的革命性变革, NLP 的前景无疑令人激动。我希望通过本教材, 能够为读者提供以语言模型为主线的 NLP 知识体系, 并让读者能深入理解大语言模型前沿理论, 掌握大语言模型实践技能。希望这本教材成为您学习与实践 NLP 的得力工具, 激发您对自然语言处理无尽的好奇心和创造力。

在本书即将付梓之际，我深感荣幸与感激，借此机会向所有在本书编写过程中给予帮助和支持的个人与机构致以最诚挚的谢意！

谨向自然语言处理领域的诸多前辈和专家致以崇高的敬意。我要特别感谢大模型技术相关领域的研究学者，正是你们在 Transformer、GPT、GLM 等模型的发展过程中所做出的杰出贡献，为本书的内容奠定了坚实的基础。你们的无私分享和合作精神，推动了这一领域的飞速发展，并为本书提供了丰富的理论依据和实践经验。

感谢中国人工智能学会自然语言理解专委会主任王小捷教授对本书编写给予的大力支持，并提出了若干宝贵的建议。感谢我们通用人工智能团队以及金融智能与金融工程四川省重点实验室的同仁们，为本书的完成提供了不可或缺的支持。本书初稿完成以后，我们团队自然语言处理课程组老师做了大量的工作。陈珍珠博士校对了第一章的部分内容。陈中普博士和陈星延博士分别校对并修改了第二章的部分内容。陈星延博士校对了第九章全部内容。陈中普博士和张阳博士分别校对了第十章的部分内容。陈中普博士、陈珍珠博士和张阳博士分别编写并校对了第十一章的部分内容。陈星延博士校对了第十二章的全部内容。陈珍珠博士校对了第十三章的全部内容。陈中普博士编写并校对了第十四章和第十五章的部分内容。施龙博士和陈星延博士分别校对了附录 A 预备知识的部分内容。此外，李庆博士和谢志龙博士也参与了本书的部分校对工作。另外，特别感谢韦鳗珍、白芊芊、罗灵、王瑞、刘银峰、许雯婷、钟一、顾添承、刘雅玲、黄浩南、唐川清、郭宇、杨闻博和邓黄怡等同学在文稿编辑、图表绘制和审稿校对等方面所付出的努力，这是本书完成的基础。

本书引用了一些优秀参考文献中的图表、公式和案例等，征求了相关作者的意见并得到了积极支持，在此表示由衷的感谢！

此外，由衷感谢机械工业出版社的辛勤付出，感谢你们在本书编写和出版过程中所展现的专业精神和不懈努力，让本书最终得以面世。

本书的编写得到了国家自然科学基金项目的资助。

最后，对所有期待本书的读者表示感谢，你们的期望和支持激励我不断深入研究。希望本书能为你们带来启发和帮助，并在你们的学术和实践道路上提供有益的参考。

再次感谢所有为本书付出心血和智慧的朋友们！愿本书能为自然语言处理领域的发展贡献绵薄之力，并激励更多的研究者投身于这一充满前景的研究领域！

由于编者水平有限，书中难免有疏漏和不足之处，恳请读者批评指正！如果您发现书中的任何错误或遇到任何问题，可以发送邮件至编者邮箱，真诚期

待您的反馈。

赵 宇
2024 年 6 月

目 录

序 言	iii
前 言	v
第一章 绪论	1
1.1 自然语言处理概述	1
1.2 自然语言处理简史	2
1.3 自然语言处理传统研究内容	4
1.3.1 传统基础技术	5
1.3.2 实际应用	31
1.4 自然语言处理与大模型发展现状	41
1.5 本书内容安排	42
1.6 讨论	43
1.7 习题	43
第一部分 语言模型基础	45
第二章 词向量	47
2.1 概述	47
2.2 独热表示	47
2.3 Word2Vec 模型	49
2.3.1 CBOW 模型	50
2.3.2 Skip-gram 模型	53
2.4 Glove 模型	54
2.5 ELMo 模型	57
2.5.1 双向语言模型	58

2.5.2	ELMo 模型	59
2.6	讨论	60
2.7	习题	60
第三章	统计语言模型	61
3.1	概述	61
3.2	统计语言模型	62
3.3	平滑技术	64
3.4	讨论	67
3.5	习题	67
第四章	神经语言模型	69
4.1	概述	69
4.2	神经概率语言模型	69
4.3	循环神经网络	72
4.3.1	循环神经网络结构	72
4.3.2	RNN 的主要结构形式	74
4.4	基于循环神经网络的语言模型	77
4.5	讨论	81
4.6	习题	81
第五章	预训练语言模型	83
5.1	概述	83
5.2	Seq2Seq 模型	83
5.2.1	模型结构	83
5.2.2	模型训练与使用技巧	87
5.3	Attention 机制	88
5.4	Transformer 模型	90
5.4.1	模型整体结构	91
5.4.2	模型推理过程	93
5.5	预训练语言模型	103
5.5.1	BERT	103
5.5.2	GPT-1	110
5.6	语言模型使用范式	113
5.6.1	预训练-微调范式	113

5.6.2	大模型-提示工程范式	113
5.7	讨论	114
5.8	习题	114
 第二部分 大模型理论		 115
第六章 大语言模型架构		117
6.1	概述	117
6.2	基于 Transformer 的模型架构	118
6.2.1	编码预训练语言模型	118
6.2.2	解码预训练语言模型	119
6.2.3	编解码预训练语言模型	132
6.3	非 Transformer 的模型架构	134
6.3.1	FAT 模型	134
6.3.2	AFT 模型	137
6.3.3	RWKV 模型	138
6.4	大模型架构配置	145
6.5	讨论	146
6.6	习题	147
 第七章 多模态大模型架构		 149
7.1	概述	149
7.2	Vision Transformer	149
7.3	CLIP	151
7.3.1	模型架构	151
7.3.2	CLIP 模型零样本分类	152
7.3.3	CLIP 其他应用	153
7.4	BLIP	153
7.4.1	BLIP 概要	154
7.4.2	BLIP 模型详解	156
7.5	BLIP2	159
7.5.1	BLIP2 概要	159
7.5.2	BLIP-2 架构	160
7.6	讨论	164
7.7	习题	165

第八章 大模型预训练	167
8.1 概述	167
8.2 预训练数据工程	167
8.2.1 数据源	168
8.2.2 多模态数据集	170
8.2.3 数据处理	174
8.2.4 预训练数据与大模型的性能关系	176
8.3 大模型预训练	177
8.3.1 预训练任务	177
8.3.2 预训练方法	182
8.3.3 优化参数设置	183
8.3.4 可扩展训练技术	184
8.4 讨论	186
8.5 习题	187
第九章 大模型微调	189
9.1 指令微调	189
9.1.1 指令微调概念	189
9.1.2 构造指令实例	190
9.1.3 多模态指令微调	192
9.1.4 指令微调优化方法	194
9.1.5 指令微调的效果	195
9.2 对齐微调	196
9.2.1 RLHF	196
9.2.2 RLHF 的发展历程	198
9.2.3 对齐微调技术	199
9.3 微调算法	204
9.4 讨论	206
9.5 习题	206
第十章 提示工程	209
10.1 概述	209
10.2 提示工程基础	209
10.2.1 提示词的组成	210
10.2.2 提示工程方法	211

10.2.3 图片提示	216
10.3 情景学习	219
10.3.1 定义	219
10.3.2 示例设计方法	220
10.3.3 有效性解释	222
10.4 思维链	223
10.4.1 提示方法	224
10.4.2 过程优化	226
10.4.3 外部引擎	229
10.5 提示工程进阶	231
10.5.1 提示破解	231
10.5.2 提示可靠性	234
10.6 讨论	235
10.7 习题	235
第十一章 涌现	237
11.1 概述	237
11.2 涌现现象	238
11.2.1 涌现的概念定义和特征	238
11.2.2 涌现的普适模型	242
11.3 大语言模型中的涌现	244
11.3.1 大语言模型中涌现的定义	245
11.3.2 大语言模型的涌现能力	246
11.3.3 大语言模型涌现能力的来源	248
11.4 缩放法则	249
11.4.1 模型性能的影响因素	250
11.4.2 缩放法则的量子化假设	250
11.5 大语言模型的可解释性	251
11.6 讨论	253
11.7 习题	253
第十二章 大模型评估	255
12.1 概述	255
12.2 评估方式	255
12.2.1 人工评估	255

12.2.2 自动评估	256
12.3 评估任务	261
12.3.1 基本评估任务	261
12.3.2 高级评估任务	265
12.3.3 评估数据集	266
12.4 评估指标	269
12.4.1 准确性	270
12.4.2 安全性	273
12.4.3 鲁棒性	274
12.4.4 高效性	277
12.4.5 其他评估指标	279
12.5 讨论	279
12.6 习题	282
第十三章 探讨	285
13.1 概述	285
13.2 基于大模型的智能体和具身智能	285
13.2.1 智能体	285
13.2.2 具身智能	287
13.3 大模型垂直领域应用	288
13.3.1 金融	288
13.3.2 法律	290
13.3.3 医疗	292
13.3.4 旅游	296
13.4 大模型的挑战与局限	298
13.4.1 幻觉现象	298
13.4.2 计算成本高昂	300
13.4.3 时效性差	301
13.4.4 专业领域表现欠佳	302
13.4.5 输出不稳定	302
13.5 大模型的社会影响	303
13.5.1 虚构事实	303
13.5.2 毒性与偏见	305
13.5.3 学术造假	306
13.5.4 环境成本	307

13.5.5 主流霸权	308
13.6 讨论	308
13.7 习题	309

第三部分 大模型实践 311

第十四章 大模型本地开发 313

14.1 概述	313
14.2 Transformers 编程基础	314
14.2.1 Transformers 关键组件	314
14.2.2 对话模型实战	316
14.3 大模型微调	318
14.3.1 使用 Transformers 微调预训练模型	319
14.3.2 使用 Transformers 微调大模型	321
14.3.3 使用 LLaMA-Factory 微调大模型	324
14.4 讨论	326
14.5 习题	326

第十五章 基于大模型的应用开发 327

15.1 概述	327
15.2 基于 OpenAI 的应用开发	327
15.2.1 关键概念	328
15.2.2 入门程序	329
15.2.3 OpenAI 模型	331
15.2.4 开发指南	334
15.2.5 应用案例	347
15.2.6 使用 Azure OpenAI	350
15.3 基于通义千问的应用开发	353
15.3.1 入门程序	353
15.3.2 通义千问模型	354
15.4 基于 LangChain 的应用开发	357
15.4.1 LangChain 入门程序	358
15.4.2 LangChain 的模型 IO	361
15.4.3 LangChain 的数据连接	365
15.4.4 LangChain 的链	368

15.4.5	LangChain 的记忆	369
15.4.6	LangChain 的代理	372
15.4.7	LangChain 的回调	372
15.4.8	LangChain 应用案例	375
15.5	讨论	379
15.6	习题	380
附录 A 预备知识		383
A.1	概率论基本概念	383
A.1.1	概述	383
A.1.2	概率	383
A.1.3	条件概率	384
A.1.4	贝叶斯法则	385
A.1.5	随机变量	386
A.1.6	二项式分布	387
A.1.7	联合概率分布和条件概率分布	387
A.1.8	期望与方差	388
A.1.9	贝叶斯决策理论	388
A.2	信息论基本概念	389
A.2.1	概述	389
A.2.2	熵	389
A.2.3	联合熵和条件熵	390
A.2.4	互信息	391
A.2.5	相对熵	393
A.2.6	交叉熵	393
A.2.7	困惑度	394
A.3	机器学习基本概念	395
A.3.1	概述	395
A.3.2	训练方式	395
A.3.3	常用算法和模型	398
A.4	强化学习基本概念	405
A.4.1	概述	405
A.4.2	强化学习中的马尔可夫过程决策	405
A.4.3	策略迭代	409
A.4.4	重要性采样	412

目 录	xvii
A.4.5 近端策略优化算法	414
附录 B 缩略语表	421
附录 C 翻译对照表	423
附录 D 相关学术会议与学术组织	429
索 引	431
参考文献	431

第一章 绪论

“自然语言处理被誉为人工智能皇冠上的明珠。”

——无名氏

1.1 自然语言处理概述

人类的语言能力是在早期儿童时期发展起来的，并在一生中不断进化。与人类不同，机器无法自然地掌握人类语言理解能力，它们需要通过如人工智能算法或者语言模型的方式，提升其在语言方面的智能。为了探究人类语言的奥秘，研究者们创立了如自然语言处理 (Natural Language Processing, NLP)、自然语言理解 (Natural Language Understanding, NLU) 和计算语言学 (Computational Linguistics, CL) 等学科，这些学科致力于使用计算方法来分析、理解和生成人类语言，从而构建人机交流的桥梁。随着时间的推移，这些研究领域在技术上取得了显著的进步，推动了人工智能在语言理解和生成方面的技术进步。

NLP、NLU 和 CL 虽然密切相关，但侧重点有所不同。NLP 是一门人工智能领域的交叉学科，涉及计算机科学、人工智能和语言学等多个领域，旨在使计算机能够解释、生成和处理人类自然语言¹。自然语言处理技术的研究不仅提高了信息处理的效率，而且也极大地拓宽了人机交互的方式，促进了科技创新和社会进步。随着技术的不断发展，NLP 在信息检索、智能问答系统、机器翻译等众多应用领域发挥着日益重要的作用。NLU 旨在使计算机能够理解和解释人类自然语言的含义和意图，其目标是让计算机能够像人类一样理解和推断自然语言文本或语音输入的意义。CL 则更侧重于语言本身的科学研究，试图通过计算机模拟来理解语言的本质和结构，包括语言的生成和理解机制。总

¹ 自然语言包括如中文、英文、俄文、阿拉伯文、西班牙语等。不同于自然语言，人们常常还使用各种编程语言等机器语言操纵计算机、使用 SQL 语言操纵数据库，以及使用手势等肢体语言沟通信息。

之，NLP、NLU和 CL都是探索人类语言奥秘、实现人机交流的重要学科，它们共同推动着语言技术的发展和應用。随着技术的不断进步，这些领域之间的界限越来越模糊。在本书中，若无特别说明，将不再区分这三个术语。

语言模型 (Language Models, LMs) 在推动 NLP、NLU和 CL的发展过程中起到了核心作用。语言模型通过计算一系列词汇在特定上下文中出现的概率，来预测下一个词的可能性。这不仅是实现机器翻译、语音识别等应用的基础，也是让机器能够更自然地与人类进行交流的关键。随着深度学习技术的发展，基于神经网络的语言模型（如 Transformers [168]）极大地推动了这些领域的发展，使得机器在理解和生成自然语言方面的能力得到了质的飞跃。除了技术的进步，数据的积累也对这些领域的发展至关重要。大数据时代的来临为这些领域的研究提供了前所未有的资源。通过分析和处理大量的文本数据，研究者们能够训练更加强大和精准的模型，这些模型能够更好地理解语言的复杂性和多样性。然而，尽管取得了巨大进步，人类语言的复杂性和深层的语义理解仍然是机器面临的巨大挑战。语言不仅仅是词汇和语法的组合，它还包含了文化、情感和隐喻等多层次的含义。因此，未来的研究将需要更深入地探索人类语言的这些方面，以及如何让机器更加准确地理解和生成具有深层含义的语言。

1.2 自然语言处理简史

自然语言处理技术的发展可以追溯到 20 世纪 50 年代初，经历了多个阶段和重要的突破。以下是自然语言处理技术发展的简单概述：

1. 早期阶段 (1950s - 1960s)

- 1950 年代早期，阿兰·图灵 (Alan Turing) 提出了著名的“图灵测试 (Turing Test)”，探讨了机器是否能够模仿人类的自然语言交流。
- 1954 年，Georgetown 大学和 IBM 合作，开发了一个名为“Georgetown-IBM 实验室翻译器”的系统，用 IBM 701 电脑将 60 多个俄语句子翻译成英文，标志着机器翻译的起步。

除了图灵测试和 Georgetown-IBM 实验之外，这个时期还诞生了一些初步尝试来理解语言的技术，如知名语言学家诺姆·乔姆斯基 (Noam Chomsky) 的生成语法理论，它对后来的语言模型和理解有着深远的影响。