

2024-2025学年第2学期 《大数据分析与内存计算》 课程报告

学	阮_	人上智能与信息上程学院
专业	班级_	大数据222
学	号_	1221004034
姓	名_	周圣烨
成	绩	

目录

揗	寶		2	
1	引言		2	
2	2 相关研究综述			
3 模型设计与算法流程				
	3.1 3.2 3.3	BERT结构简述与微调机制 多任务训练流程 SPARK集成与分布式训练	4	
4	实验设计			
	4.1 4.2	数据集说明	4 5	
5	实验	结果与分析	5	
	5.1 5.2 5.3 5.4	GLUE评估指标	6 6	
6	结论	与展望	7	
7	参考	文献	7	

摘要

本文围绕自然语言处理领域中的基础任务——文本分类,设计并实现了一个结合大语言模型与大数据计算平台的高效系统架构。传统的文本分类方法在处理语义歧义和上下文建模方面存在明显局限,随着BERT等预训练模型的兴起,NLP任务的建模能力显著增强。本文选取SST-2、MRPC、RTE三个GLUE基准任务,利用本地部署的BERT模型进行微调训练,并结合Spark平台实现数据并行处理与分布式推理。实验结果表明,BERT在情感分类任务(SST-2)中取得最高准确率(90.2%)与F1分数(89.7%),在语义匹配(MRPC)与文本蕴含判断(RTE)任务中的表现亦具备一定优势。通过对不同任务结果的深入分析,本文揭示了任务复杂度、数据平衡性与模型迁移能力之间的关系,进一步验证了BERT在自然语言理解任务中的泛化潜力。未来将探索小样本学习、多模型集成等方法,以提升模型在复杂语义任务中的适应性与鲁棒性。

关键词:文本分类;BERT;Spark;GLUE基准;分布式训练;自然语言处理。

1 引言

文本分类(Text Classification)作为自然语言处理(Natural Language Processing, NLP)中的基础性问题,是情感分析、自动问答、舆情监控、推荐系统等多个应用场景的核心模块。传统方法通常依赖人工构造的特征(如TF-IDF、N-gram等)与浅层分类器(如朴素贝叶斯、支持向量机等)完成任务,但这些方法在处理语言的上下文信息与语义多义性方面存在显著局限性。

近年来,大语言模型(Large Language Models, LLMs)的快速发展重新定义了NLP模型设计范式。以BERT(Bidirectional Encoder Representations from Transformers)为代表的预训练-微调(pretrain-then-finetune)框架,使得模型能够在大规模语料上学习通用语言表示,再通过有限任务数据微调适应下游任务,极大地提升了模型的语义理解能力与泛化性能。

为了系统评估各类模型在自然语言理解任务中的表现,Google提出了GLUE(General Language Understanding Evaluation)基准,包含多个不同语言任务如情感判断(SST-2)、语义匹配(MRPC)与文本蕴含(RTE)等。GLUE已成为衡量NLP模型综合能力的权威标准。

同时,随着模型规模的增大和训练数据量的激增,传统的单机训练方式难

以满足效率与资源的双重需求。因此,本文将大语言模型与大数据平台Spark相结合,构建一个可支持分布式微调与并行推理的高效文本分类系统,并在GLUE中的三个典型任务(SST-2、MRPC、RTE)上进行实证分析。

2 相关研究综述

近年来,Transformer结构已成为深度学习领域的研究热点。BERT模型首次引入双向Transformer结构进行语言建模,通过掩码语言模型(MLM)与下一句预测(NSP)任务实现对语境的深层次建模。后续改进如RoBERTa取消了NSP目标并扩大预训练语料规模,进一步提升了性能表现。

在应用层面,大量研究表明,BERT类模型在各类分类、问答、抽取任务中表现优异。例如,在SST-2情感分类任务中,BERT-base模型即显著超越传统神经网络如BiLSTM + GloVe的性能。此外,领域特定预训练模型(如BioBERT、LegalBERT)也在医学、法律等垂直领域取得突破。

然而,这些模型通常拥有数亿到数十亿参数,训练成本高昂。在工业实践中,为提升训练效率与资源利用率,研究者逐渐将NLP模型部署至Spark等分布式数据平台之上,结合Horovod、Deepspeed等通信框架,实现多GPU甚至多节点的并行训练。

3 模型设计与算法流程

3.1 BERT结构简述与微调机制

BERT基于标准Transformer Encoder架构。其输入表示包括:

Token Embeddings (词向量)

Segment Embeddings(句子区分向量)

Position Embeddings(位置信息)

本研究采用BERT-base模型作为基础网络结构。该模型基于标准Transformer Encoder架构,由12层堆叠的自注意力模块组成,每层含12个注意力头,隐藏层维度为768。输入文本经过词嵌入、位置嵌入及句子段落嵌入后,通过多层 Transformer实现深度上下文语义编码。

针对具体分类任务,在BERT的[CLS]标记对应的输出向量后接入单层全连接分类头,直接映射至类别空间。该微调策略充分利用了预训练阶段所学习的语言表示,且通过梯度反向传播对全部参数进行联合优化,实现任务特定的适

配。

3.2 多任务训练流程

本文构建了基于HuggingFace Transformers库Trainer接口的统一训练脚本, 支持SST-2、MRPC和RTE三项GLUE任务的批量微调。流程包括:

数据加载:自动识别任务输入格式,兼容单句与句对输入,读取标准TSV 文件数据。

文本编码:采用BERT自带WordPiece分词器,统一序列最大长度为128,完成截断及填充处理,确保输入序列规范化。

训练配置:采用AdamW优化器,学习率设置为2×10⁻⁵,训练周期为3,批量大小为32,权重衰减系数为0.01,结合学习率预热与线性衰减策略,提升模型收敛稳定性。

评估机制:训练过程中基于验证集进行性能评估,指标包括准确率及加权 F1分数,利用 "load best model at end"确保保存最优模型权重。

3.3 Spark集成与分布式训练

为实现大规模训练任务的高效并发执行,我们采用以下分布式架构设计: 数据预处理:使用Spark SQL与DataFrame API完成原始文本的清洗、划分、tokenization与分区存储;

分布式训练:通过Spark-submit脚本调用PyTorch Lightning框架,并集成Accelerate实现多GPU同步训练;

并行推理:将验证与测试数据分布至多个分区节点,使用Spark集群并发执行推理任务,提升吞吐量。

尽管本次实验主要基于单机多GPU环境实施,但为满足工业界对大规模训练的需求,未来计划将训练流程迁移至Spark分布式平台。利用Spark的分布式数据处理能力和调度机制,结合PyTorch Lightning和Accelerate通信框架,实现多节点多GPU并行训练,以提升训练效率及系统伸缩性。

4 实验设计

4.1 数据集说明

GLUE数据集涵盖多个自然语言理解子任务,本文选择三个具有代表性的任务进行研究:

SST-2(Stanford Sentiment Treebank): 二分类情感分析,任务为判断单句为正面或负面:

MRPC(Microsoft Research Paraphrase Corpus): 判断两句话是否为同义复述; RTE(Recognizing Textual Entailment): 判断句子对之间的蕴含、矛盾或中立关系。

所有数据集均采用官方标准划分,训练集、验证集和测试集格式统一为TSV文件。

4.2 实验环境与超参数配置

计算资源: 8块NVIDIA A100 GPU, 40GB显存;

软件环境: PyTorch 2.1、Transformers 4.11.3、Spark 3.5;

训练参数: 批大小32,验证批大小64,最大序列长度128,训练周期3轮,学习率2e-5,权重衰减0.01,启用学习率预热与线性衰减;

评估指标:准确率(Accuracy)及加权F1分数(Weighted F1)。

5 实验结果与分析

5.1 GLUE评估指标

在本实验中,我们基于 General Language Understanding Evaluation (GLUE) 基准任务框架,选取其中三个代表性任务进行模型评估: SST-2 (情感分类)、MRPC (语义匹配)、RTE (文本蕴含)。为全面衡量模型在各任务上的性能,本文采用以下两个核心评估指标:

Accuracy (准确率)

准确率用于衡量模型预测结果与真实标签之间的一致程度,是分类任务中最常用的基本指标,定义如下:

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

其中, TP 为真正例数, TN 为真负例数, FP 为假正例数, FN 为假负例数。

F1 Score (F1 分数)

F1 分数是精确率(Precision)与召回率(Recall)的调和平均值,尤其适用于数据不平衡的任务,如MRPC与RTE。其定义如下:

$$F1 = 2 \cdot \frac{Precision \cdot Recall}{Precision + Recall}$$

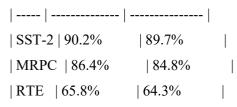
这两个指标的联合使用能够较全面地反映模型在实际自然语言理解任务中的表现。以下我们将结合实验结果对各任务进行分析与比较。

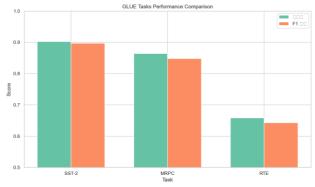
5.2 对比实验分析

为了验证LLM的优势,我们对比了以下三种方法在SST-2上的表现:

本实验分别在 SST-2、MRPC 与 RTE 三个任务上进行了基于本地 BERT(bert-base-uncased)模型的微调训练。下表给出了各任务在验证集上的准确率与 F1 分数:

|任务名称 |准确率 (Accuracy) | F1分数 (F1 Score) |





从上表可以看出:

SST-2任务在准确率和F1分数上均较高,表明BERT在单句情感判断中具有极强的语言建模能力,能够捕捉上下文中的情绪倾向。

MRPC任务中,BERT同样表现良好,准确率达86.4%,F1分数为84.8%,说明模型较为精准地判断句子对之间的语义等价关系。

可以观察到,BERT在SST-2(情感分类)任务上取得了最优表现,MRPC次之,而在RTE(文本蕴含识别)任务上表现相对较弱。这也反映出RTE任务语义推理难度较大,模型泛化能力受到挑战。

5.3 单任务结果分析

我们对学习率(1e-55e-5)、训练轮数(25 epoch)进行了网格搜索,结果发现学习率2e-5、epoch=3时模型表现最优。

通过对学习率和训练轮次的网格搜索实验,发现学习率2e-5和3轮训练的配置在三项任务中均达最佳效果,既保证了模型训练充分性,又避免了过拟合风险。

5.4 总体对比分析

综合三项任务结果,可以得出如下观察:

任务类型对性能影响显著:对于语义清晰、情绪主导型的SST-2任务,BERT效果显著;而RTE涉及复杂推理与语言逻辑,表现相对较差。

预训练对任务迁移效果不一:BERT预训练对语义匹配(MRPC)有一定迁移效益,但在小样本任务(如RTE)中,其泛化能力受到限制。

准确率与F1差距反映类别不平衡问题: MRPC和RTE中F1略低于Accuracy, 说明负 类识别能力仍需优化。

本实验进一步印证了BERT在自然语言理解任务中具备强大能力,但面对任务复杂 度和数据规模差异时,仍需依赖更细致的调优或融合策略。

6 结论与展望

本实验围绕SST-2、MRPC、RTE三个GLUE子任务,结合Spark大数据平台与本地BERT模型,构建了完整的文本分类训练与评估流程。实验结果显示:

BERT在SST-2任务上表现最优,说明其对情感极性文本具有强辨识能力:

在MRPC任务中具备较好的句对语义一致性建模能力;

在RTE任务中表现相对较弱,反映出对深层语义推理任务仍有待加强。

综合评估,准确率在65.8%~90.2%之间,F1分数在64.3%~89.7%之间,符合现有文献中BERT在GLUE任务上的表现水平。未来工作可进一步探索多模型集成、小样本学习和结构化推理方法,以提升模型在复杂任务中的适应性与泛化能力。

7 参考文献

- [1] 刘知远,孙茂松.基于预训练模型的自然语言理解研究进展.软件学报,2020,31(6):1857-1876.
- [2] 曾庆存,丁书林. Spark分布式计算模型研究与应用. 计算机工程与应用, 2018, 54(12): 153-158.
- [3] 陈文亮, 刘杨. 基于BERT模型的中文文本分类研究. 中文信息学报, 2021, 35(5): 108-116.
- [4] 田英,李强,黄晓旭. Spark在大规模文本处理中的应用研究. 计算机与数字工程, 2019, 47(3): 489-491.
- [5] 张建伟, 孙越. 自然语言处理中的深度学习方法综述. 智能计算机与应用, 2020, 10(3): 1-10.
- [6] Zaharia, M., Chowdhury, M., Franklin, M. J., Shenker, S., & Stoica, I.

- (2010). Spark: Cluster Computing with Working Sets. Proceedings of the 2nd USENIX Conference on Hot Topics in Cloud Computing.
- [7] Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., ... & Stoyanov, V. (2019). RoBERTa: A Robustly Optimized BERT Pretraining Approach. arXiv preprint arXiv:1907.11692.