

填表说明

一、本报告中相关的技术或数据如涉及知识产权保护、军工项目保密等内容，请作脱密处理。

二、请用宋体小四字号撰写本报告，可另行附页或增加页数，A4纸双面打印。

三、表中所涉及的签名都必须用蓝、黑色墨水笔，亲笔签名或签字章，不可以打印代替。

四、同行专家业内评价意见书编号由工程师学院填写，编号规则为：年份4位+申报工程师职称专业类别(领域)4位+流水号3位，共11位。

一、个人申报

(一) 基本情况【围绕《浙江工程师学院（浙江大学工程师学院）工程类专业学位研究生工程师职称评审参考指标》，结合该专业类别(领域)工程师职称评审相关标准，举例说明】

1. 对本专业基础理论知识和专业技术知识掌握情况(不少于200字)

在研究生期间，我系统地学习了计算机技术领域和数据科学领域的基础理论知识和专业技术，尤其在深度学习、机器学习、数据挖掘、图像处理、大语言模型、强化学习等前沿方向积累了扎实的理论功底与丰富的实践经验。我熟练掌握了深度学习框架PyTorch，能够开发和训练各种类型的神经网络模型，如卷积神经网络（CNN）、图卷积神经网络(GCN)和Transformer，尤其在医学影像数据分析和多模态数据融合方面有深入的实践经验。此外，我关注大语言模型的最新发展，深入研究了基于人类反馈的强化学习（RLHF），并掌握了如何运用AI自动标注策略提高标注效率和质量。我还对特征交互学习进行深入研究，提出并开发了创新的旋转因子分解机算法，有效解决了高阶特征交互的计算瓶颈，研究成果已被国际顶级会议KDD2024所认可。通过多个科研与工程项目的参与，我进一步巩固并拓展了对专业知识的理解与应用，建立了较为完善的专业技术知识体系，提升了分析真实问题并提出创新技术解决方案的能力，为从事复杂系统研发和工程创新奠定了坚实基础。

2. 工程实践的经历(不少于200字)

在研究生阶段，我参与了多个工程实践项目，积累了丰富的项目实战经验。其中，在“无创预测弥漫大B细胞淋巴瘤（DLBCL）基因突变状态”项目中，我作为算法组成员，参与数据预处理、模型构建与训练等核心任务。该项目旨在利用深度学习技术，通过影像和临床数据预测患者基因突变状态，为精准医疗提供辅助决策工具。

此外，在“RLAIF采样检测优化系统”项目中，我作为主要开发者之一，设计并实现了基于采样思想的AI标注数据筛选与修正模块。通过数据一致率评估与样本质量优化，大幅降低了人工标注成本，提高了系统在实际大模型强化学习微调中的可用性。该成果已申请软件著作权，并作为论文的工程实现支撑。

我还参与了特征交互学习项目，开发了旋转因子分解机（RFM），有效解决了高阶特征交互的计算出现指数爆炸以及缺乏特征之间相依性学习的问题，成功在国际会议KDD2024发表论文，充分展现了自身的工程实践能力与创新潜力。

通过这些实践，我不仅提升了编程与建模能力，还增强了在复杂系统中解决实际问题、跨专业协作和成果转化的能力，形成了良好的工程意识和项目管理能力。

3. 在实际工作中综合运用所学知识解决复杂工程问题的案例（不少于1000字）

在实际工作中，我始终致力于将计算机专业的理论知识与工程实践相结合，解决复杂、动态的实际问题。以下两个项目案例，展现了我在 AI 系统开发与优化中的综合应用能力和工程创新思维。高质量训练数据一直都是人工智能的基石，这两个案例就聚焦于提高数据质量。

案例一：RLAIF 采样检测优化系统 —— 构建 AI 标注数据的质量防线

问题背景与挑战：大规模语言模型技术的成功在很大程度上依赖于强化学习微调，高质量的标注数据是强化学习微调的关键。由于人工标注成本高昂，学术界提出了利用现有 AI 模型对数据标注以降低成本方案。对于高质量的 AI 标注数据，直接采用可以显著降低成本并保持后续模型训练性能。然而，AI 标注数据的质量在不同任务中可能存在显著差异，直接使用这些数据可能导致模型性能下降甚至引发安全隐患。面对这一挑战，我主导设计并开发了“RLAIF

采样检测优化系统”，提出了一种通用的采样框架，通过抽样检测和修正策略，针对不同标注质量的数据提供有效的处理方案。

技术方案：该框架首先随机抽取少量的 AI 标注数据重新进行人工标注，以此估计 AI 标注与人类偏好的一致性水平，并根据一致率将 AI 标注数据分为低、中、高三种情况，分别采取不同的处理策略：在低一致率时拒绝使用 AI 标注数据，避免低质量数据对模型训练产生负面影响；在高一致率时直接使用 AI 标注数据进行强化学习微调，抽样检测过程进一步提升了 AI 标注数据的可信性；在中等一致率时，通过采样修正部分数据样本，显著提升数据质量，使其达到高质量数据的标准，从而实现对 AI 标注样本的有效利用。我从理论上证明了抽样检测估计能够有效反映 AI 标注数据的总体一致率，并分析了抽样检测过程的渐进和非渐进性质。此外，我们还推导了在中等数据质量情况下，样本修正阶段所需修正样本数量的相合性，为算法的实际应用提供了理论支持。

实践应用：对于实际应用与测试，我们选择四个任务作为测试样本：字符串反转仿真实验、IMDB 负面影评生成、有害性回答生成和有用性回答生成。在每个任务中，我们将未经处理的 AI 标注数据与处理后的数据分别用于强化学习微调，测试 LLM

性能表现。实验结果表明：

高一一致性数据可直接用于训练，效果接近人工标注；中一致性数据经过采样修正后，性能几乎等同于人工数据，而系统整体将人工标注成本降低了超过 60%。

成果转化与意义：该项目已成功落地为“RLAIF 采样检测优化系统 V1.0”，获得软件著作权。

系统为 AI 数据增强与安全强化学习提供了工程化工具支持，推动了数据自动化处理在工程实践中的可行性与高效性，在理论与实践均具有突破性意义。

案例二：Rotative Factorization Machines —— 破解特征交互的表达瓶颈（KDD 2024）

问题背景与挑战：在推荐系统中，特征交互建模是构建高质量输入特征、进而提升模型预测准确率的关键步骤。在推荐系统、广告点击率（CTR）预测等核心任务中，模型的输入往往由用户属性、内容特征、上下文环境等多源异构信息构成。这些特征在原始状态下通常是稀疏的、离散的，并不足以直接反映用户偏好与行为之间的复杂关系。因此，对原始特征进行交互建模，生成更加语义丰富、表达力更强的组合特征，是实现高质量推荐的关键步骤。然而，现有主流方法如 FM、DeepFM 等大多只建模二阶特征交互，缺乏对高阶组合的表达能力。而尝试引入高阶交互的模型又面临参数维度指数级增长的问题，极易导致计算瓶颈和训练不稳定，难以在大规模推荐系统中落地应用。为了解决这一长期存在的技术难题，我们提出了 Rotative Factorization Machines (RFM)，并将该成果发表于国际顶级会议 KDD 2024。

技术方案：RFM 的核心思想是将特征映射到复数空间，通过角度编码 ($e^{i\theta}$) 实现交互权重的旋转表达。模型将输入特征编码为 d 维向量并转换为角度形式，再通过多个自注意力旋转模块建模特征间高阶依赖关系。最后，利用欧拉公式将复数空间映射回实数域，并通过模学习模块捕捉交互强度。我们还从理论上证明了该方法在保持高表达能力的同时，能有效避免特征阶数指数增长的问题。

实践应用：在五个真实数据集上开展了大规模实验，涵盖了点击率预测、电影推荐与应用推荐等典型场景：

Criteo：广告点击率预测任务，包含超过 4500 万条样本与 39 个字段；

Avazu：广告点击率预测竞赛数据；

MovieLens-1M 与 ML-Tag：电影评分与标签任务；

Frappe：基于上下文的应用推荐数据集。

实验表明，Rotative Factorization Machines (RFM) 在 AUC 和 LogLoss 指标上全面超越主流方法，不仅显著提升了高阶特征交互建模能力，同时通过线性梯度增长设计，有效避免了传统方法中的计算爆炸问题，兼顾了精度、稳定性与效率。

成果转化与意义：项目成果以论文形式被 KDD 2024 接收，并在巴塞罗那做了口头汇报，

获得关注。RFM 方法为高阶特征交互建模提供了新的思路，也推动了“低计算复杂度+高表达能力”在推荐系统中的融合应用。

总结：理论支撑与工程落地并重，解决真实 AI 系统难题

围绕“数据可靠性”与“特征交互学习”两大挑战，我系统地完成了从理论建模、算法设计到工

程实现的探索与创新。RLAIF 项目解决了 AI 自动标注数据质量不稳定的问题，为大模型强化学习微调提供了安全、高效、低成本的数据支撑；RFM 项目突破了高阶特征交互建模的表达瓶颈，显著提升了推荐系统中特征交互后的数据质量与模型性能。

通过这两个项目，我不仅深化了对强化学习、特征建模、复数编码、自注意力机制等关键理论的理解，还锤炼了从问题抽象、技术实现到工程落地的完整实践能力，展现了扎实的理论功底、系统设计思维与创新解决问题的能力，为在复杂 AI 工程场景中持续推动技术应用与创新打下了坚实基础。

(二) 取得的业绩(代表作)【限填3项, 须提交证明原件(包括发表的论文、出版的著作、专利证书、获奖证书、科技项目立项文件或合同、企业证明等)供核实, 并提供复印件一份】

1. 公开成果代表作【论文发表、专利成果、软件著作权、标准规范与行业工法制定、著作编写、科技成果获奖、学位论文等】

成果名称	成果类别 [含论文、授权专利(含发明专利申请)、软件著作权、标准、工法、著作、获奖、学位论文等]	发表时间/授权或申请时间等	刊物名称/专利授权或申请号等	本人排名/总人数	备注
RLAIF采样检测优化系统 V1.0	计算机软件著作权	2025年03月03日	登记号: 2025SR0369 356	2/2	
中国研究生数学建模竞赛	获奖	2023年01月01日			

2. 其他代表作【主持或参与的课题研究项目、科技成果应用转化推广、企业技术难题解决方案、自主研发设计的产品或样机、技术报告、设计图纸、软课题研究报告、可行性研究报告、规划设计方案、施工或调试报告、工程实验、技术培训教材、推动行业发展中发挥的作用及取得的经济社会效益等】

(三) 在校期间课程、专业实践训练及学位论文相关情况	
课程成绩情况	按课程学分核算的平均成绩： 88 分
专业实践训练时间及考核情况(具有三年及以上工作经历的不作要求)	累计时间： 1 年(要求1年及以上) 考核成绩： 88 分
本人承诺	
<p>个人声明：本人上述所填资料均为真实有效，如有虚假，愿承担一切责任，特此声明！</p> <p style="text-align: right;">申报人签名：石雨鸿</p>	

浙江大学研究生院
攻读硕士学位研究生成绩单

学号: 22260290	姓名: 石雨鸿	性别: 女	学院: 工程师学院	专业: 计算机技术	学制: 2.5年						
毕业时最低应获: 24.0学分		已获得: 26.0学分		入学年月: 2022-09	毕业年月:						
学位证书号:			毕业证书号:			授予学位:					
学习时间	课程名称	备注	学分	成绩	课程性质	学习时间	课程名称	备注	学分	成绩	课程性质
2022-2023学年秋季学期	新时代中国特色社会主义思想理论与实践		2.0	92	公共学位课	2022-2023学年冬季学期	计算机视觉		2.0	72	跨专业课
2022-2023学年秋季学期	工程技术创新前沿		1.5	89	专业学位课	2022-2023学年春季学期	数学建模		2.0	86	专业选修课
2022-2023学年秋季学期	数据科学技术与软件实现		2.0	92	专业学位课	2022-2023学年春夏学期	工程伦理		2.0	97	公共学位课
2022-2023学年秋冬学期	研究生论文写作指导		1.0	90	专业学位课	2022-2023学年夏季学期	自然辩证法概论		1.0	85	公共学位课
2022-2023学年秋冬学期	数据分析的概率统计基础		3.0	96	专业选修课	2023-2024学年夏季学期	研究生英语		2.0	免修	公共学位课
2022-2023学年秋冬学期	高阶工程认知实践		3.0	90	专业学位课	2023-2024学年夏季学期	研究生英语基础技能		1.0	免修	公共学位课
2022-2023学年冬季学期	产业技术发展前沿		1.5	90	专业学位课		硕士生读书报告		2.0	通过	

说明: 1. 研究生课程按三种方法计分: 百分制, 两级制 (通过、不通过), 五级制 (优、良、中、及格、不及格)。
2. 备注中“*”表示重修课程。

学院成绩校核章:

成绩校核人: 张梦依

打印日期: 2025-06-03

成绩校核章

中华人民共和国国家版权局 计算机软件著作权登记证书

证书号： 软著登字第15025554号

软件名称： RLAIIF采样检测优化系统
V1.0

著作权人： 浙江大学

权利取得方式： 原始取得

权利范围： 全部权利

登记号： 2025SR0369356

根据《计算机软件保护条例》和《计算机软件著作权登记办法》的规定，经中国版权保护中心审核，对以上事项予以登记。



2025年03月03日

软件著作权情况说明

软件名称: RLAIIF 采样检测优化系统 V1.0

著作权人: 浙江大学

现就本软件的开发者排序、申请人贡献情况及与论文的相关性说明如下:

一、开发者排序及贡献情况

本软件由以下人员共同开发完成 (按贡献排序):

- 梁克维: 提供算法建模与理论分析的指导, 参与系统关键模块的思路制定, 对软件实现提供重要支持。
- 石雨鸿: 负提出系统总体设计思路, 主导抽样检测算法与采样修正策略的核心模块实现, 负责数据实验的设计与结果分析, 完成主要开发工作及系统整合。

二、申请人贡献说明

申请人石雨鸿在本软件的开发过程中发挥了主导作用, 提出了主要的系统设计思路, 独立完成了系统的核心功能开发, 并主导了多轮实验验证工作, 对软件的形成做出了实质性贡献。

三、与论文的相关性说明

本软件为论文《一种基于采样的强化学习微调大模型的方法》中提出方法的具体实现系统, 论文主要研究了在利用 AI 自动标注数据进行强化学习微调的背景下, 如何通过**抽样检测和样本修正策略**提升 AI 标注数据质量, 降低人工标注成本。

论文提出的三种数据一致率处理策略 (拒绝使用、直接使用、采样修正) 在本系统中均有完整实现。本软件通过以下几个模块体现了论文的理论成果:

- 一致率估计模块:** 随机抽取样本进行人工标注对比, 估算总体一致率;
- 分级处理策略模块:** 根据一致率判断, 分别采取拒绝、使用或修正处理;
- 采样修正模块:** 在中等一致率下对部分数据样本进行人工修正, 提升整体质量;
- 性能验证模块:** 集成半仿真与真实数据集评估方法, 实现理论与实验的闭环。

该软件为论文研究成果的具体应用载体, 为验证论文中提出的方法提供了有力支撑。

四、导师签署意见

本软件开发工作由我校研究生完成, 成果已在相关论文中体现, 符合学校著作权归属规范。

导师签字: 梁克维

日期: 2025 年 3 月 22 日



中国研究生创新实践系列大赛

“中国光谷·华为杯”第十九届中国研究生数学建模竞赛

“Optics Valley Of China·Huawei Cup” The 19th China Post-Graduate Mathematical Contest in Modeling

获奖证书

浙江大学 石雨鸿

荣获“中国光谷·华为杯”第十九届中国研究生数学建模竞赛

三等奖



主办单位：中国学位与研究生教育学会



中国科协青少年科技中心



中国研究生数学建模竞赛组委会



承办单位：浙江大学

编号：F2022300974

二〇二三年一月