

1 测试与评估

本章旨在评估项目所开发的伪造新闻检测系统，特别是引入核心的因果干预算法后的性能表现。

1.1 测试环境与数据

- **训练数据集：** 使用了包含 23975 条新闻样本的训练数据集进行训练。该数据由平台提供。

- **测试数据集：** 使用了包含 499 条新闻样本的测试数据集进行评估。该数据集由平台提供。

- 文本检测训练数据使用流程：

- 1 • 数据预处理阶段开始于原始数据清洗，运行 `data_converter.py` 脚本移除 HTML 标签和特殊字符，统一编码格式，处理缺失值，标准化文本，合并标题与正文，并进行初步标签转换。随后使用将清洗后的数据按 8:2 比例分割为训练集和验证集。

- 2 • 特征增强处理分为两大步骤：

(1) RAG 知识库构建和数据标注。知识库构建通过 `build_knowledge_base.py` 收集权威事实核查数据、已验证新闻、谣言模式库和权威信源列表，并创建向量数据库索引。数据标注使用 `label_data.py` 提取三类关键特征：内容特征（关键词、实体、主题、事实性）、风格特征（情感、语言风格、修辞手法、专业度）和传播特征（时间信息、来源可靠性、传播模式、社会影响力）。

(2) Chain of Thought 推理增强包括两个环节：通过 `generate_cot_templates.py` 构建标准化推理模板，涵盖信息来源分析、事实核查步骤、逻辑推理链路和可信度评估；然后使用 `generate_reasoning.py` 为每条训练数据生成详细推理过程，包含多步骤推理链、证据支持、逻辑关系和结论推导。

- 3 • 因果干预训练分三步实现：

(1) 用 `_build_causal_graph` 函数构建特征间因果关系网络；

(2) 通过 `_generate_content_counterfactuals` 函数、`_generate_style_counterfactuals` 函数和 `_generate_propagation_counterfactuals` 函数分别生成内容、风格和传播维度的反事实样本；

(3) 最后使用 `_calculate_causal_effects` 函数计算各特征的直接、间接、总体和交互因果效应。

- 4 • 模型训练阶段使用 `train_model.py` 对已经预下载好的模型进行多阶段训练，包括基础特征、CoT 增强、因果干预优化和模型集成，随后用 `evaluator.py` 以适用于新闻真假检测场景的二分类评估方法进行全

面评估。最后进行结果验证与优化，通过 `test_model.py` 生成性能报告。这一完整流程确保模型充分利用 RAG 知识检索、CoT 思维链推理和因果干预分析等先进技术，实现高精度假新闻检测能力。

- 多模态检测训练数据使用流程
- 使用 DGM4 团队上传在 Huggingface 上的 20 万行数据集进行训练，并保存训练日志

1.2 评估指标

主要采用以下指标来衡量模型性能：

- **准确率 (Accuracy)**：模型正确预测（真或假）的样本占总样本的比例。
- **精确率 (Precision)**：模型预测为“伪造”的样本中，实际确实为“伪造”的比例。
- **召回率 (Recall)**：实际为“伪造”的样本中，被模型成功检测出来的比例。
- **F1 分数 (F1-Score)**：精确率和召回率的调和平均值，综合反映模型的稳健性。
- **运行时间**：包括总运行时间和平均每样本处理时间，评估模型的效率。
- **样本分布**：观察模型预测结果中真/假新闻的分布情况。

1.3 测试方法

采用了对比实验的方法，比较了两种算法版本在同一测试数据集上的表现：

1. **基础版**：未使用因果干预算法的模型（仅使用 cot 和 rag 增强的模型）。
2. **因果干预增强版**：引入了因果干预分析机制的模型。

通过对比两者的性能指标，评估因果干预算法对检测性能的实际提升效果。

1.4 测试结果

测试结果主要如下：

指标	基础版（无因果干预）	因果干预增强版	提升幅度
准确率 (Accuracy)	79.96%	95.12%	+15.16 百分点
精确率 (Precision)	71.98%	96.23%	+24.25 百分点
召回率 (Recall)	97.99%	97.08%	-0.91 百分点

F1 分数 (F1-Score)	82.99%	96.65%	+13.66 百分 点
预测为假新闻	339 条 (67.94%)	324 条 (64.93%)	
预测为真新闻	160 条 (32.06%)	175 条 (35.07%)	
总运行时间	79.45 秒	92.36 秒	+12.91 秒
平均每样本处理时 间	159.22 毫秒	185.09 毫秒	+25.87 毫秒

1.5 结果分析

从测试结果对比来看，引入因果干预算法带来了显著的性能提升：

- **准确率大幅提升：**准确率从约 80% 提升至 95% 以上，提升了 15.16 个百分点，表明因果干预显著增强了模型区分真伪新闻的整体能力。
- **精确率质的飞跃：**精确率提升了惊人的 24.25 个百分点，达到 96.23%。这意味着模型在判断为“伪造”时非常可靠，极大地减少了将真实新闻误判为伪造的情况（False Positives）。这对于维护平台或用户的信任至关重要。
- **召回率保持高位：**召回率略有下降（-0.91 百分点），但仍保持在 97% 以上的极高水平，说明模型检测出真正伪造新闻的能力依然非常强，几乎没有漏判。
- **F1 分数显著提高：**F1 分数提升了 13.66 个百分点，达到 96.65%，表明模型在精确率和召回率之间取得了更好的平衡，整体性能更加稳健。
- **预测分布更均衡：**因果干预版预测的真假新闻比例（约 65% vs 35%）比基础版（约 68% vs 32%）略显均衡，结合精确率的大幅提升，表明模型减少了“过度警惕”将真实新闻误判为伪造的倾向。
- **效率代价可接受：**虽然引入因果干预使得平均处理时间增加了约 16.25%（从 159ms 增加到 185ms），但考虑到准确率提升了 15.16%，精确率提升了 24.25%，性能上的巨大收益远超时间成本的增加。对于许多应用场景，这种程度的效率牺牲是完全可以接受的。

1.6 综合评估

测试结果有力地证明了**因果干预算法**在本项目的伪造新闻检测任务中的有效性。它不仅显著提高了检测的准确性和可靠性（尤其是在减少误报方面），而且在保持高召回率的同时，大幅提升了整体 F1 分数。虽然处理时间略有增加，但性能提升带来的价值远超其成本。这表明基于因果干预的方法能够更好地捕捉新闻真伪的本质特征，而不仅仅是表面关联，从而构建出更强大、更可信的检测系统。