

## 1. 实验概览

- **模型架构**：基于BERT的 `bert_base_uncased` 模型，采用本地部署（`local_modal=true`）。
- **数据集**：GLUE基准中的CoLA（Corpus of Linguistic Acceptability）数据集，任务类型为单句分类（`TaskForSingleSentenceClassification`）。
- **训练参数**：训练轮次（`epochs=3`），随机种子（`random_seed=42`），启用GPU加速（`use_gpu=true`）。未启用任何防御机制（`defenderEnabled=false`）。
- **执行攻击类型**：
  - **安全攻击**：数据投毒攻击（`PoisoningAttack`）。
  - **隐私攻击**：无执行记录（模型反演、梯度反演、模型窃取攻击均为空）。

## 2. 基准性能分析

- **正常训练结果**：`normalTrain` 字段为空列表，无基准性能数据可用。无法提供准确率（accuracy）和F1分数分析。
- **与文献比较**：BERT-base在CoLA数据集的典型基准准确率约80-85%（参考GLUE排行榜）。实验中缺失基准数据，需假设攻击前的模型性能符合文献范围，以评估攻击影响。

## 3. 安全攻击分析

- **数据投毒攻击(PoisoningAttack)**：
  - **攻击解释**：通过污染训练数据集（如注入误导性样本），降低模型泛化能力。本次攻击参数：投毒率15%（`poisoning_rate=0.15`），训练轮次15（`epochs=15`）。
  - **结果分析**：
    - 攻击后任务准确率=40%，F1分数=37.5%（`resultData=[ "0.4", "0.375" ]`）。
    - **性能下降显著**：对比BERT在CoLA的典型基准（~80%），攻击导致准确率下降约50%。F1分数（37.5%）表明模型在精确率和召回率上均严重退化。
    - **趋势分析**：仅单次攻击记录，无法计算多次攻击的变化趋势。但低值（<40%）表明高攻击有效性。
  - **平均性能下降**：因无基准数据，假设攻击前性能为80%，则平均下降幅度达50%。

## 4. 隐私攻击分析

- **模型反演(RLMI)、梯度反演(FET)、模型窃取(ModelStealingAttack)**：无执行记录，跳过分析。

## 5. 横向对比分析

### 安全攻击特性对比

评估维度	投毒
性能影响	高（准确率下降50%）
检测难度	中等（需数据审计）
缓解成本	高（需重训模型）

### 对比分析：

- 数据投毒攻击直接导致模型性能崩溃（准确率40%），**破坏强度远超常规安全威胁**。
- 隐蔽性中等：投毒率15%属常见阈值，但未触发防御机制（`defenderEnabled=false`）。

### 隐私攻击特性对比

无隐私攻击数据，跳过表格与分析。

### 关键风险评估

1. **最大业务威胁**：数据投毒攻击（导致模型失效，业务中断风险）
2. **最高合规风险**：训练数据完整性丧失（违反数据治理规范）
3. **最紧急漏洞**：训练管道污染（投毒率15%即可生效）

## 6. 防御建议

- **针对性防御措施**：
  - **数据投毒攻击**：部署训练数据清洗（如离群值检测）、差分隐私（添加噪声）、或对抗训练（增强鲁棒性）。启用防御后需复测投毒率容忍阈值。
- **监控指标**：
  - 实时跟踪训练集分布偏移（如KL散度）。
  - 部署模型性能降级警报（如准确率单次下降>20%时触发）。
- **架构改进**：
  - 采用联邦学习分散数据风险。
  - 集成模型验证模块（如运行时输入过滤）。

## 7. 结论

本次实验聚焦数据投毒攻击的安全影响，使用BERT-base模型在CoLA数据集上执行攻击。关键发现如下：

1. **攻击有效性极高**：投毒率15%时，模型准确率骤降至40%，F1分数仅37.5%，表明攻击成功破坏模型核心功能。结合无防御机制（defenderEnabled=false），此漏洞构成最高业务风险，需优先修补。
  2. **基准数据缺失的局限**：因normalTrain为空，无法量化性能下降绝对值，需假设基准准确率为80%。建议后续实验强制包含正常训练以建立可靠基线。
  3. **防御紧迫性**：数据投毒攻击易于实施且成本低（仅需污染训练集），但缓解需重训模型，成本高昂。推荐部署实时数据审计工具（如CleanLab）以降低风险。
  4. **实验扩展方向**：
    - 测试不同投毒率（如5%-30%）对性能的影响曲线。
    - 启用防御策略（如差分隐私）并评估其有效性。
    - 补充其他攻击类型（如对抗样本）以全面评估模型弱点。
- 未来研究**应探索自适应投毒攻击（如动态调整污染样本）及跨模型迁移攻击的鲁棒性。

### 附录：指标解释

- **准确率(Accuracy)\*\***：分类任务中正确预测的比例，计算公式为  $(TP+TN)/(TP+TN+FP+FN)$ 。
- **F1分数**：精度(Precision)和召回率(Recall)的调和平均，综合评估模型性能，计算公式为  $2 \times (Precision \times Recall) / (Precision + Recall)$ 。
- **投毒率(Poisoning Rate)\*\***：污染样本占训练集的比例，直接影响攻击强度。