

1. 实验概览

- **模型架构:** BERT-Base-Uncased (本地部署)
- **数据集:** IMDB影评数据集 (单句分类任务)
- **训练参数:** CPU训练, 随机种子42, 未启用GPU加速
- **执行攻击类型:**
 - 安全攻击: 对抗样本攻击 (AdvAttack)、后门攻击 (BackdoorAttack)、数据投毒攻击 (PoisoningAttack)
 - 隐私攻击: 模型反演攻击 (RLMI)、梯度反演攻击 (FET)、模型窃取攻击 (ModelStealingAttack)

2. 基准性能分析

- **正常训练结果:** 准确率90.0%, F1分数0.859
- **基准对比:** BERT在IMDB分类任务文献基准为89-92%准确率, **当前模型性能符合预期**, 未出现显著偏差。

3. 安全攻击分析

• 对抗样本攻击(AdvAttack)

通过同义词替换和字符干扰欺骗模型。关键发现:

- 平均攻击成功率65.0% (两次攻击均成功)
- **攻击后准确率暴跌65%** (100%→35%/98%→35%)
- 攻击分布: 平均成功12.5次/失败7.5次/跳过0次
- **异常:** 启用防御后攻击成功率未降低, **防御策略可能失效**

• 后门攻击(BackdoorAttack)

植入触发词改变模型行为。关键发现:

- 毒化后准确率异常降至0.277% (原始0.883%)
- 隐蔽性指标全为 NaN, **数据严重缺失导致无法评估**
- **需核查攻击是否实际执行** (executed: false)

• 数据投毒攻击(PoisoningAttack)

污染训练数据以降低模型性能。关键发现:

- 攻击后准确率降至50.0% (-40%), F1降至0.759
- **性能下降幅度达44.4%** (对比基准90.0%)
- 投毒率15%时即造成显著破坏

4. 隐私攻击分析

• 模型反演攻击(RLMI)

通过模型输出反推训练数据。关键发现:

- 攻击/推理阶段成功率均达100%
- 词错误率高 (攻击阶段0.794, 推理阶段0.720)
- **成功率与错误率呈负相关**, 揭示输出质量与隐私泄露的权衡

• 梯度反演攻击(FET)

从梯度重建原始数据。关键发现:

- 最终指标: ROUGE-1(0.728)/ROUGE-2(0.0)/ROUGE-L(0.437)

- **词汇恢复率0%，编辑距离85（极高），完全恢复率0%
- 训练动态：ROUGE-L从0.437骤降至0.0（Epoch 2）
- 关键转折点：**首Epoch后分数崩塌，**攻击完全失败**

• 模型窃取攻击(ModelStealingAttack)

通过查询复制目标模型。关键发现（启用防御 vs 未启用）：

- 受害者模型准确率：0.8862 vs 0.8000
- 窃取模型准确率：0.5020（两者相同）
- **模型相似度仅52.84%**，未达窃取阈值
- 训练过程：损失从0.48降至0.14，但验证准确率稳定在90%
- 防御未显著提升安全性**（相似度无变化）

5. 横向对比分析

安全攻击特性对比

评估维度	对抗样本	后门	投毒
性能影响	高	中	高
检测难度	中	低	高
缓解成本	高	中	低

对比分析：

- 对抗样本攻击成功率**65%，破坏强度超投毒攻击20%
- 后门攻击因数据缺失无法评估隐蔽性

隐私攻击特性对比

评估维度	模型反演	梯度反演	模型窃取
信息质量	中	低	低
实施复杂度	低	高	高
防御可行性	高	中	中

对比分析：

- 模型反演成功率**100%，但信息质量受限
- 梯度反演完全恢复率**0%，技术可行性存疑

关键风险评估

- 最大业务威胁：**对抗样本攻击（直接导致模型失效）
- 最高合规风险：**模型反演攻击（100%数据重建成功率）
- 最紧急漏洞：**数据投毒（44.4%性能下降且易实施）

6. 防御建议

- **对抗样本攻击:**
 - 部署对抗训练 (Adversarial Training)
 - 监控指标: **输入扰动敏感度**、置信度波动
- **后门攻击:**
 - 采用输入过滤 (如ONION) 清除触发词
 - 架构改进: **神经元激活分析**检测异常模式
- **数据投毒:**
 - 强化数据清洗流程
 - 实时监控: **训练损失异常突增**
- **隐私攻击通用防御:**
 - 梯度裁剪+差分隐私
 - 限制模型查询频率

核心结论:

1. **对抗样本攻击构成最高业务风险** (65%成功率)
2. 模型反演攻击暴露严重隐私漏洞 (100%成功率)
3. 现有防御对模型窃取和对抗样本效果不足

附录: 指标解释

- **准确率(Accuracy)**: 模型正确预测的比例
- **F1分数**: 精确率与召回率的调和平均
- **攻击成功率**: 成功欺骗模型的攻击样本占比
- **词错误率(WER)**: 重建文本与原文本的差异率
- **ROUGE**: 文本重建质量评估 (*ROUGE-L*关注最长公共子序列)
- **编辑距离**: 两字符串间的最小编辑操作次数
- **模型相似度(Agreement)**: 窃取模型与原始模型预测一致性