

1. 实验概览

- 模型架构：BERT-base-uncased
 - 数据集：IMDb（电影评论情感分类）
 - 训练参数：CPU训练、随机种子42、单句分类任务
 - 攻击类型及分类：
 - 安全攻击：对抗样本攻击（AdvAttack）、后门攻击（BackdoorAttack）、投毒攻击（PoisoningAttack）
 - 隐私攻击：模型反演（RLMI）、梯度反演（FET）、模型窃取（ModelStealingAttack）
-

2. 基准性能分析

- 正常训练结果：
 - 准确率：90.000%
 - F1分数：0.859
 - 基准对比：在IMDb数据集上，BERT-base通常达到94%准确率（文献基准），当前模型性能偏低4%，可能存在优化不足或数据偏差。
-

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。

- 平均攻击成功率：65.0%（两次攻击均成功）
- 攻击后准确率暴跌：启用防御时从100.0%→35.0%；未启用时从98.0%→35.0%
- 攻击分布：共25次尝试（成功20次，失败15次，跳过0次）
- 关键发现：防御机制未显著降低攻击成功率（均为65%），需强化文本鲁棒性。

后门攻击(BackdoorAttack)：

植入触发词操纵模型输出。

- 毒化前后准确率：0.883% → 0.277%（暴跌68.6%）
- 隐蔽性评估：困惑度(PPL)、语义相似性(USE)、语法正确性均为NaN（数据异常）
- 毒化效果雷达图概念：

```
radarChart
  title 后门攻击隐蔽性
  "准确率下降" : 68.6
  "困惑度" : 0
  "语义相似性" : 0
  "语法正确性" : 0
```

注：非数值指标因数据缺失无法展示

数据投毒攻击(PoisoningAttack)：

污染训练数据以降低模型性能。

- 攻击效果：准确率从90.000%→50.000% (-40.0%)，F1从0.859→0.759 (-11.6%)
 - 性能下降幅度：平均降低40%准确率，投毒率15%即造成严重破坏。
-

4. 隐私攻击分析

模型反演攻击(RLMI)：

重构训练数据隐私内容。

- 攻击阶段成功率100.00%，词错误率0.79411
- 推理阶段成功率100.00%，词错误率0.71951
- 平衡分析：高成功率但高词错误率，表明重建文本可识别主题但细节失真。

梯度反演攻击(FET)：

通过梯度泄露原始数据。

- 最终指标：

ROUGE-1 ROUGE-2 ROUGE-L 词汇恢复率 编辑距离 完全恢复率

0.728 0.0 0.437 1.0 85 0%

- 训练动态：ROUGE-1在首epoch即达峰值(0.728)，后续无提升
- 关键转折点：首epoch后分数停滞，完全恢复率为0%暴露攻击局限性。

模型窃取攻击(ModelStealingAttack)：

复制受害者模型功能。

- 性能对比：

模型类型 启用防御 准确率 相似度

受害者模型 是 88.62% -

受害者模型 否 80.00% -

窃取模型 - 50.20% 52.84%

- 训练过程：损失从0.4808→0.1456，训练准确率从86.25%→92.71%，但验证准确率稳定在90.00%
 - 核心风险：窃取模型性能不足受害者模型的60%，相似度仅52.84%，复制效果有限。
-

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	极高 (-65%)	极高 (-68.6%)	高 (-40%)
检测难度	中 (依赖扰动)	低 (触发明显)	高 (潜伏期长)
缓解成本	中 (需对抗训练)	高 (需数据清洗)	低 (输入过滤)

对比分析：

- 对抗样本攻击成功率65%，破坏性强于其他攻击
- 后门攻击隐蔽性数据缺失，需补充评估

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	低 (高错误率)	中 (部分恢复)	低 (低相似度)
实施复杂度	高	极高	中
防御可行性	中 (输出模糊化)	低 (梯度压缩)	高 (查询限制)

对比分析：

- 梯度反演完全恢复率0%，实际泄露风险可控
- 模型窃取导致知识产权风险值达52.84%，需优先防御

关键风险评估

1. 最大业务威胁：对抗样本攻击（实时欺骗风险）
 2. 最高合规风险：模型窃取（知识产权泄露）
 3. 最紧急漏洞：投毒攻击（训练数据污染）
-

6. 防御建议

- 对抗样本：
 - 部署对抗训练（如PGD）
 - 监控指标：输入文本扰动敏感度、置信度波动
 - 后门攻击：
 - 采用异常检测（如ONION）
 - 监控指标：触发词激活频率、输出分布偏移
 - 投毒攻击：
 - 强化数据清洗（如K近邻去噪）
 - 监控指标：训练损失突变、验证集性能断层
 - 隐私攻击通用防御：
 - 梯度裁剪 + 差分隐私（针对梯度反演）
 - API查询频率限制（针对模型窃取）
 - 架构改进：
 - 集成多任务学习提升鲁棒性
 - 添加置信度校准层过滤异常输出
-

附录：指标解释

- ROUGE：衡量生成文本与参考文本的n-gram重叠率（1/L分别表示单字/双字/最长公共子序列）
- 编辑距离：两文本间最小编辑操作数（插入、删除、替换），值越大差异越大
- 词错误率(WER)：语音识别指标借用于文本重建，计算错误词占比
- 相似度(agreement)：窃取模型与受害者模型预测一致性的比例
- 困惑度(PPL)：衡量文本概率分布不确定性，值越低越“自然”