

1. 实验概览

- 模型架构: BERT-Base-Uncased
- 数据集: IMDB (情感分析任务)
- 训练参数:
 - 随机种子: 42
 - 未启用GPU加速
 - 本地模型及数据集

- 攻击类型及分类:

攻击类型	分类
AdvAttack (对抗样本)	安全攻击
BackdoorAttack (后门)	安全攻击
PoisoningAttack (投毒)	安全攻击
RLMI (模型反演)	隐私攻击
FET (梯度反演)	隐私攻击
ModelStealingAttack (窃取)	隐私攻击

2. 基准性能分析

- 正常训练结果:
 - 准确率: **90.0%**
 - F1分数: **0.859**
- 与基准对比:

BERT在IMDB的典型基准准确率约92-95% (如Devlin et al., 2019)，本实验90.0%略低，可能因训练配置 (如未用GPU) 或数据局部性导致。

3. 安全攻击分析

对抗样本攻击 (AdvAttack)

通过替换同义词/干扰字符欺骗模型 (策略: *TextFooler*)。

- 平均攻击成功率: 65.0% (两次攻击均值)
- 准确率变化:
 - 启用防御时: 100.0% → 35.0% (**下降65%**)
 - 未启用防御时: 98.0% → 35.0% (**下降63%**)
- 攻击分布:

状态	次数
成功	25
失败	15

跳过	0
关键发现: 防御启用对成功率无显著影响, 需检查防御策略有效性。	

后门攻击 (BackdoorAttack)

植入触发词改变模型行为 (策略: *TrojanLM*) 。

- **毒化前后准确率:**

◦ 原始数据集: 0.883% → 毒化后: 0.277% (**下降68.6%**)

- **隐蔽性评估:**

- 困惑度 (PPL) : NaN (数据异常)
- 语义相似性 (USE) : NaN
- 语法正确性: NaN

异常标注: 关键指标缺失, 可能因攻击未执行 (executed: false) 。

数据投毒攻击 (PoisoningAttack)

污染训练数据以降低性能 (投毒率: 15%) 。

- **性能下降:**

- 准确率: 90.0% → 50.0% (**下降44.4%**)
- F1分数: 0.859 → 0.759 (**下降11.6%**)

- **平均性能下降:** 准确率 **-44.4%**, F1 -11.6%

推论: 投毒对准确率破坏更显著, 符合预期。

4. 隐私攻击分析

模型反演攻击 (RLMI)

通过模型输出反演原始数据 (策略: 强化学习优化) 。

- **阶段对比:**

阶段	成功率	词错误率
攻击阶段	100.00%	0.79411
推理阶段	100.00%	0.71951

- **平衡分析:**

成功率100%但词错误率高 (>0.7), 表明恢复数据**可用性低**。

梯度反演攻击 (FET)

从梯度反演输入文本 (策略: *FET/SWAT*) 。

- **最终综合指标:**

指标	值
ROUGE-1	0.7275

ROUGE-2	0.0
ROUGE-L	0.4365
词汇恢复率	1.0
编辑距离	85
完全恢复率	0%

- 训练动态：

- ROUGE-1峰值0.7275（首Epoch），后续无提升。

关键转折点：首Epoch后分数停滞，反演未收敛。

模型窃取攻击 (ModelStealingAttack)

通过查询窃取模型 (策略：MeaeQ)。

- 模型性能对比：

模型	准确率
受害者模型	88.62%
窃取模型	50.20%

- 相似度：Agreement = **0.5284** (低相似性)

- 训练过程：

- 损失从0.4808降至0.1456 (10轮迭代)。
- 训练准确率从86.25%升至92.71%，但验证集稳定在90.0%。

发现：窃取模型严重过拟合，泛化能力差。

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 (↓65%)	高 (↓68.6%)	中 (↓44.4%)
检测难度	中	高	低
缓解成本	低	高	中

对比分析：

- 对抗样本攻击成功率**65%，破坏性最强；
- 后门攻击因指标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	低 (词错误率高)	中 (ROUGE-1:0.72)	低 (窃取准确率50%)

实施复杂度	高	极高	中
防御可行性	中	低	高

对比分析：

- 梯度反演完全恢复率0%，但词汇恢复率100%，存在部分敏感信息泄露风险；
- 模型窃取导致知识产权风险（相似度52.84%）。

关键风险评估

1. **最大业务威胁**: 对抗样本攻击 (直接破坏服务可用性)
2. **最高合规风险**: 梯度反演 (潜在训练数据泄露)
3. **最紧急漏洞**: 数据投毒 (易实施且性能下降显著)

6. 防御建议

- **对抗样本**:
 - 防御: 对抗训练 (如PGD)
 - 监控: 输入文本扰动检测 (如词频突变)
- **后门攻击**:
 - 防御: 输入过滤 (如ONION策略)
 - 监控: 触发词异常激活率
- **投毒攻击**:
 - 防御: 数据清洗 (基于离群值检测)
 - 监控: 训练集分布偏移指标
- **模型反演**:
 - 防御: 限制模型输出置信度
 - 架构: 差分隐私层
- **梯度反演**:
 - 防御: 梯度压缩/添加噪声
 - 监控: 梯度幅值异常
- **模型窃取**:
 - 防御: API查询限流/水印技术
 - 架构: 模型混淆 (Obfuscation)

附录：指标解释

- **准确率 (Accuracy)** : 正确预测样本比例。
- **F1分数**: 精确率与召回率的调和平均。
- **攻击成功率**: 成功欺骗模型的样本占比。
- **词错误率 (WER)** : 反演文本与原文本的字错误率。
- **ROUGE**: 文本生成质量评估 (ROUGE-1/2/L关注n-gram匹配)。
- **编辑距离**: 两文本间最小编辑操作次数。
- **Agreement**: 窃取模型与受害者模型预测一致性。