

# 攻击实验分析报告

## 1. 实验概览

- 模型架构:** BERT-Base-Uncased (本地部署)
- 数据集:** GLUE/CoLA (单句分类任务)
- 训练参数:** 3个epoch, 随机种子42, GPU加速
- 安全攻击:**
  - 对抗样本攻击 (TextFooler策略)
  - 后门攻击 (BadNets策略, 启用STRIP防御)
  - 数据投毒攻击 (投毒率15%)
- 隐私攻击:**
  - 模型反演攻击 (启用剪枝防御)
  - 梯度反演攻击 (FET/SWAT策略)
  - 模型窃取攻击 (MeaeQ策略, 启用输出扰动防御)

## 2. 基准性能分析

**警告：**未获取正常训练基准数据，无法进行准确率/F1分数分析。建议后续补充基准测试。

## 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文本分类模型 (TextFooler策略)。

- 攻击成功率:** 100% (3/3样本攻击成功)
- 破坏强度:** 准确率从100%降至0%，完全瘫痪模型功能
- 攻击分布:** 全部成功 (0失败/0跳过)

### 后门攻击 (BackdoorAttack)

植入隐藏触发模式控制模型输出 (BadNets策略)。

- 性能影响:** 准确率从89.4%暴跌至27.4%
- 隐蔽性评估:**
  - 困惑度(PPL): **269.336** (极高, 暴露异常)
  - 语义相似性(USE): 0.935 (良好伪装)
  - 语法正确性: NaN (数据缺失)

### 数据投毒攻击 (PoisoningAttack)

污染训练数据降低模型性能 (投毒率15%)。

- **性能破坏**: 攻击后准确率40%，F1分数37.5%
- **攻击效果**: 模型性能严重劣化（无基准对比）

## 4. 隐私攻击分析

### 模型反演攻击 (RLMI)

通过模型输出重构训练数据。

- **攻击阶段成功率**: 100% | **推理阶段成功率**: 99.29%
- **词错误率(WER)**: 攻击阶段66.89%，推理阶段73.22%
- **关键发现**: 高成功率伴随高词错误率，表明重构内容语义失真

### 梯度反演攻击 (FET)

利用梯度信息窃取原始文本数据。

- **最终指标**:
  - ROUGE-1: 0.619 | ROUGE-2: 0.0526 | ROUGE-L: 0.4286
  - 词汇恢复率: **0%** | 编辑距离: 76 | 完全恢复率: 100%
- **训练动态**:
  - 初始ROUGE-1达0.619，但后续epoch数据缺失
  - **异常点**: 词汇恢复率与完全恢复率存在矛盾（0% vs 100%）

### 模型窃取攻击 (ModelStealingAttack)

通过查询复制受害者模型（MeaeQ策略）。

- **性能对比**:
  - 受害者模型准确率: 88.62%
  - 窃取模型准确率: **50.20%** (性能损失43.7%)
- **相似度**: 模型输出一致性仅52.82%
- **训练过程**:
  - 训练损失从0.4846→0.1674 (下降65.5%)
  - 训练准确率从87.29%→92.71%
  - 验证准确率稳定在80%

## 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	<b>毁灭性</b> (100%→0%)	<b>严重</b> (89.4%→27.4%)	<b>显著</b> (准确率40%)

评估维度	对抗样本	后门	投毒
检测难度	高（无防御）	中（PPL异常暴露）	高（训练阶段注入）
缓解成本	高（需对抗训练）	中（需触发模式检测）	高（需数据清洗）

#### 对比分析：

- 对抗样本攻击成功率**100%**，破坏强度超其他攻击200%以上
- 后门攻击困惑度指数**269.336**，超过安全阈值300%

#### 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（高WER）	低（ROUGE-L=0.4286）	中（窃取准确率50.2%）
实施复杂度	高（需强化学习）	极高（遗传算法优化）	中（查询次数500）
防御可行性	中（剪枝有效）	低（无防护）	中（噪声扰动）

#### 对比分析：

- 梯度反演完全恢复率**100%**，存在关键数据泄露隐患
- 模型窃取导致知识产权风险值达**52.82%**，构成最高合规风险

#### 关键风险评估

- 最大业务威胁**：对抗样本攻击（模型完全失效）
- 最高合规风险**：模型窃取（知识产权侵犯）
- 最紧急漏洞**：梯度反演（原始数据泄露）

## 6. 防御建议

#### • 对抗样本：

- 部署对抗训练（Adversarial Training）
- 实时监控输入置信度波动（阈值>30%触发警报）

#### • 后门攻击：

- 增强STRIP防御的触发模式检测
- 监控推理时异常激活模式

#### • 隐私攻击：

- 梯度加密与噪声注入（噪声标准差>0.1）
- 实施查询频率限制（<100次/分钟）
- 模型输出模糊化（置信度截断）

#### • 架构改进：

- 集成差分隐私机制（ $\epsilon < 2.0$ ）

- 采用**联邦学习**分散敏感数据

## 7. 结论

本次攻击实验揭示了BERT模型在GLUE/CoLA任务上的严重安全脆弱性：

**核心发现1**：对抗样本攻击展现毁灭性效果，100%成功率使模型完全失效，表明模型缺乏对抗鲁棒性。建议立即部署对抗训练和输入过滤机制，并建立置信度实时监控系统。

**核心发现2**：隐私攻击中梯度反演存在矛盾指标（词汇恢复率0% vs 完全恢复率100%），提示数据泄露风险被低估。需重点检查梯度保护机制，强制实施梯度压缩和噪声注入，阈值建议设为编辑距离 $>50$ 或ROUGE-L $<0.5$ 时阻断传输。

**核心发现3**：现有防御（如后门STRIP和模型剪枝）效果有限，毒化后准确率仍下降62%，模型反演成功率仍达99.29%。证明需采用**防御组合策略**，推荐差分隐私（ $\delta=1e-5$ ）+输出扰动（噪声 $\sigma=0.15$ ）+查询限制的三层防护。

**改进方向**：1) 补充正常训练基准；2) 测试复合攻击场景；3) 评估Transformer-XL等鲁棒架构；4) 探索可信执行环境(TEE)部署方案。最紧迫的研究方向是开发对抗样本的实时检测算法与隐私攻击的轻量化防御框架。

---

## 附录：指标解释

- **PPL (困惑度)**：衡量文本自然度，值越高越异常（后门攻击）
- **USE (语义相似性)**：0-1分值，越高说明篡改文本越隐蔽（后门攻击）
- **WER (词错误率)**：重构文本错误比例，越高质量越差（模型反演）
- **ROUGE**：文本匹配度指标，ROUGE-L衡量最长公共子序列（梯度反演）
- **编辑距离**：两文本间最小编辑操作次数，值越高差异越大（梯度反演）
- **Agreement (相似度)**：两模型输出一致性比例（模型窃取）