

## 基于数据特征的成员推理增强攻击方法

牛俊<sup>1</sup>, 沈括<sup>2</sup>, 王月琮<sup>3</sup>, 韩雪雪<sup>3</sup>, 张嘉伟<sup>1</sup>, 李兴华<sup>1</sup>, 张玉清<sup>1,2,4</sup>

(1. 西安电子科技大学网络与信息安全学院, 陕西 西安 710126; 2. 西安电子科技大学杭州研究院, 浙江 杭州 311231;  
3. 西安电子科技大学通信工程学院, 陕西 西安 710071; 4. 国家计算机网络入侵防范中心 (中国科学院大学), 北京 101408)

**摘要:** 针对现有成员推理攻击 (MIA) 方法只在单一实例或场景下检测攻击的有效性, 导致检测结果不准确并误导攻击者判断这一问题, 首先构造了多个测试场景并联合多个评估指标, 设计了一种成员推理攻击有效性检测算法, 发现现有攻击方法的假阳率 (FPR) 最高可达 100%。在此基础上, 根据数据特征和攻击原理, 提出了一种基于数据特征的成员推理增强攻击方法, 以降低攻击的假阳率并提高攻击性能。实验结果表明, 所提方法在 4 个基准数据集上, 可将 15 个现有成员推理攻击的假阳率从 100% 降低到 0, 将精确率从 58.27% 提高到 100%, 将成员增益从 32.37% 提高到 100%, 进一步增强了攻击的有效性。

**关键词:** 机器学习; 成员推理攻击; 数据安全; 数据特征

**中图分类号:** TP391

**文献标志码:** A

**DOI:** 10.11959/j.issn.1000-436x.2025162

## Enhanced membership inference attack method leveraging data characteristics

NIU Jun<sup>1</sup>, SHEN Kuo<sup>2</sup>, WANG Yuecong<sup>3</sup>, HAN Xuexue<sup>3</sup>,  
ZHANG Jiawei<sup>1</sup>, LI Xinghua<sup>1</sup>, ZHANG Yuqing<sup>1,2,4</sup>

1. School of Cyber Engineering, Xidian University, Xi'an 710126, China

2. School of Hangzhou Institute of Technology, Xidian University, Hangzhou 311231, China

3. School of Telecommunications Engineering, Xidian University, Xi'an 710071, China

4. National Computer Network Intrusion Protection Center (University of Chinese Academy of Sciences), Beijing 101408, China

**Abstract:** Aiming at the problem that the existing membership inference attack (MIA) methods only present generic testing results and had not sufficiently evaluated the effectiveness of such attacks, which could mislead attackers, multiple test scenarios were first constructed and various evaluation metrics were combined to propose a membership inference attack validity detection algorithm. It was revealed that the false positive rate (FPR) of current attack methods could reach as high as 100%. Based on data characteristics and attack principles, an enhanced membership inference attack method was subsequently proposed to reduce the false positive rate and improve attack effectiveness. Experimental results on four benchmark datasets demonstrate that the FPR of 15 existing membership inference attacks is reduced from 100% to 0, while the attack precision is improved from 58.27% to 100%, and the membership advantage is increased from 32.37% to 100%, further enhancing the effectiveness of the attack.

**Keywords:** machine learning, membership inference attack, data security, data characteristics

收稿日期: 2025-05-13; 修回日期: 2025-09-15

通信作者: 张玉清, zhangyq@nipc.org.cn

基金项目: 国家自然科学基金资助项目 (No.U2336203, No.62125205, No.U23A20303); 国家重点研发计划基金资助项目 (No.2023YFB3106400, No.2023QY1202); 北京市自然科学基金资助项目 (No.L242015, No.4242031); 西电概念验证基金资助项目 (No.GNYZ2024QC009); 中央高校基本科研业务费专项基金资助项目 (No.ZYTS25066); 陕西省自然科学基金基础研究计划基金资助项目 (No.2024JC-YBMS-544)

**Foundation Items:** The National Natural Science Foundation of China (No.U2336203, No.62125205, No.U23A20303), The National Key Research and Development Program of China (No.2023YFB3106400, No.2023QY1202), Beijing Natural Science Foundation (No.L242015, No.4242031), The Innovation Fund of Xidian University (No.GNYZ2024QC009), The Fundamental Research Funds for the Central Universities (No.ZYTS25066), The Natural Science Basic Research Program of Shaanxi (No.2024JC-YBMS-544)

## 0 引言

算力、算法和数据的不断进步,促进了机器学习(ML, machine learning)的快速发展,进而推动了人工智能(AI, artificial intelligence)的繁荣与创新。机器学习作为AI的一种实现方式,在图像识别、自然语言处理、数据挖掘<sup>[1]</sup>、欺诈识别<sup>[2]</sup>、医学分析等领域取得了巨大进展。然而,机器学习服务(MLaaS, machine learning as a service)模型的训练需要大量数据,而这些数据涉及隐私信息,如位置轨迹、医疗和金融消费等,其严重威胁着机器学习的数据安全。同时,即使机器学习模型和算法是可信的,机器学习模型的输出(如信任分数向量)也很容易包含训练数据的隐私信息,从而导致其易遭受各种隐私攻击,如模型提取攻击、模型投毒攻击、后门攻击、对抗样本攻击、成员推理攻击(MIA, membership inference attack)<sup>[3-4]</sup>以及模型版权确权等。因此,关注和研究机器学习中的数据安全具有至关重要的意义。

成员推理攻击<sup>[3,5-7]</sup>主要推测一个数据样本是否在目标模型的训练集中,极大地威胁着训练数据的隐私安全。尽管目前已有大量成员推理攻击方法<sup>[3-7]</sup>依赖目标模型的输出概率、真实标签、影子模型和数据集的差异性等信息,来捕捉成员关系并提升攻击性能。然而,现有工作只给出了一般的测试结果,并未充分检测这些攻击方法的有效性和实用性,进而导致错误的推测结果误导攻击者并造成巨大的经济损失。为解决上述问题,本文首先根据数据特征构造了多个测试场景(TS, test scenarios),提出了一种成员推理攻击有效性检测算法,其先对目标数据进行分类并构造了4种测试场景,接着计算不同测试场景下各个攻击方法的攻击效果,并与预先设置的假阳率(FPR, false positive rate)阈值和成员增益(MA, membership advantage)阈值进行比较来判断攻击的有效性。检测出现有MIA存在较高的假阳率(如100%),而且不同数据集和类别的假阳率也不同。此外,根据攻击时使用的目标数据、数据特征和攻击原理,提出了一种基于数据特征的成员推理增强攻击方法,其主要根据攻击者能容忍的最小假阳率以及数据样本间的距离,对样本进行摘选和替换。在4个基准数据集上的实验发现,利用本文方法可极大地提升15个现有先进MIA的攻击性能。例如,将假阳率降低到0,将精

确率和成员增益等评估指标的数值提高到100%,有效增强了攻击性能。本文的主要贡献包括以下3个方面。

1) 提出了一种成员推理攻击有效性检测算法,发现现有方法存在较高的假阳率,最高可达100%,而且不同数据集和类别具有不同的假阳率。

2) 提出了一种基于数据特征的成员推理增强攻击方法,其根据目标数据特征(如模型输出概率或数据样本间的距离)及攻击原理,来挑选合适的目标数据以实施攻击,极大地增强了攻击效果。

3) 对比15种现有MIA表明,本文方法在4个基准数据集上可极大增强各个攻击的攻击效果。例如,将攻击假阳率降低到0,将精确率和成员增益提高到100%。

## 1 相关工作

### 1.1 成员推理攻击

目标ML模型 $F$ 被描述为在某个目标训练数据集 $D_{\text{train}}$ 上训练的ML模型。MIA旨在确定数据点 $x$ 是否为目标训练数据集 $D_{\text{train}}$ 中的成员。MIA包括三类:基于二元分类器的MIA利用影子技术训练一个二元攻击模型;基于评估机制的MIA根据预先定义的成员评估机制(如模型输出阈值)来实施攻击;基于数据集差异的MIA不需要训练影子模型,仅利用2个数据集间的距离差异来识别成员关系。

### 1.2 威胁模型

成员推理攻击的威胁模型如表1所示,其中, $\sqrt$ 和 $\times$ 分别指需要或不需要该知识实施攻击。由表1可知,黑盒、灰盒和白盒攻击的背景知识依次增强,攻击者的攻击能力也逐步增强。

表1 成员推理攻击的威胁模型

威胁模型	训练数据分布	目标模型		
		输入	输出	架构和参数
黑盒	$\times$	$\sqrt$	$\sqrt$	$\times$
灰盒	$\sqrt$	$\sqrt$	$\sqrt$	$\times$
白盒	$\sqrt$	$\sqrt$	$\sqrt$	$\sqrt$

### 1.3 成员推理防御

已有许多防御机制相继被提出以抵御MIA,如正则化技术<sup>[3,8]</sup>、差分隐私<sup>[9]</sup>、信任分数掩蔽技

术<sup>[10]</sup>、知识蒸馏<sup>[11-12]</sup>以及成员样本预先排除技术<sup>[13]</sup>。

### 1.4 攻击方法概述

表2展示了本文研究的15种MIA方法,其中,√和×分别指需要或不需要该知识实施攻击。BlindMI-w/o<sup>[7]</sup>和BlindMI-w攻击<sup>[7]</sup>都基于2个数据集之间的距离差异性来识别成员关系,主要区别是非成员数据集不同。GradDiff攻击<sup>[14]</sup>在联邦学习场景下,利用数据集间的距离差异性来实施攻击。NN攻击<sup>[3]</sup>利用多个影子模型的输出训练攻击模型。ImproveMI攻击<sup>[15]</sup>利用数据集不同子组的模型输出分别训练多个攻击模型。EnhancedMI攻击<sup>[16]</sup>在联邦学习场景下提出了一种增强攻击,与NN攻击类似。Top3 NN攻击<sup>[5]</sup>利用Top3输出概率来实施攻击。Top1 Threshold攻击<sup>[5]</sup>基于Top1信任分数值来鉴别成员。Loss Threshold攻击<sup>[6]</sup>比较目标样本的交叉熵损失与成员样本的平均交叉熵损失的相对大小。Label-Only攻击<sup>[6]</sup>采用了预测标签和真实标签的一致性。Top2+True攻击<sup>[7]</sup>结合了Top2输出概率和真实标签。LiRA攻击<sup>[17]</sup>联合了样本的难度校准分数和高斯似然估计来识别成员关系。Risk score攻击<sup>[18]</sup>利用样本的隐私风险分数和预测熵来实施攻击。Calibrated score攻击<sup>[19]</sup>根据样本的难度校准分数来进行识别。Distillation based攻击<sup>[20]</sup>利用蒸馏技术来设置同时依赖模型和数据的攻击阈值。

表2 成员推理攻击方法描述

攻击方法	威胁模型			影子模型	攻击模型
	黑盒	灰盒	白盒		
NN <sup>[3]</sup>	√	×	×	多个	多个
Top3 NN <sup>[5]</sup>	√	×	×	一个	一个
BlindMI-w <sup>[7]</sup>	√	√	√	×	×
BlindMI-w/o <sup>[7]</sup>	√	√	√	×	×
Label-Only <sup>[6]</sup>	√	×	×	×	×
Top2+True <sup>[7]</sup>	√	×	×	多个	多个
Top1 Threshold <sup>[5]</sup>	√	×	×	×	×
Loss Threshold <sup>[6]</sup>	√	×	×	一个	×
EnhancedMI <sup>[16]</sup>	√	×	×	多个	多个
GradDiff <sup>[14]</sup>	√	√	√	×	×
ImproveMI <sup>[15]</sup>	√	×	×	多个	多个
LiRA <sup>[17]</sup>	√	×	×	多个	×
Risk score <sup>[18]</sup>	√	×	×	多个	×
Calibrated score <sup>[19]</sup>	√	×	×	多个	×
Distillation based <sup>[20]</sup>	√	×	×	多个	×

## 2 基于数据特征的成员推理增强攻击方法

针对现有MIA并未充分检测攻击有效性,导致错误推测并造成经济损失的问题,本文首先提出了一种成员推理攻击有效性检测算法。然后,根据目标数据、数据特征和攻击原理,提出了一种基于数据特征的成员推理增强攻击方法(如图1所示)。

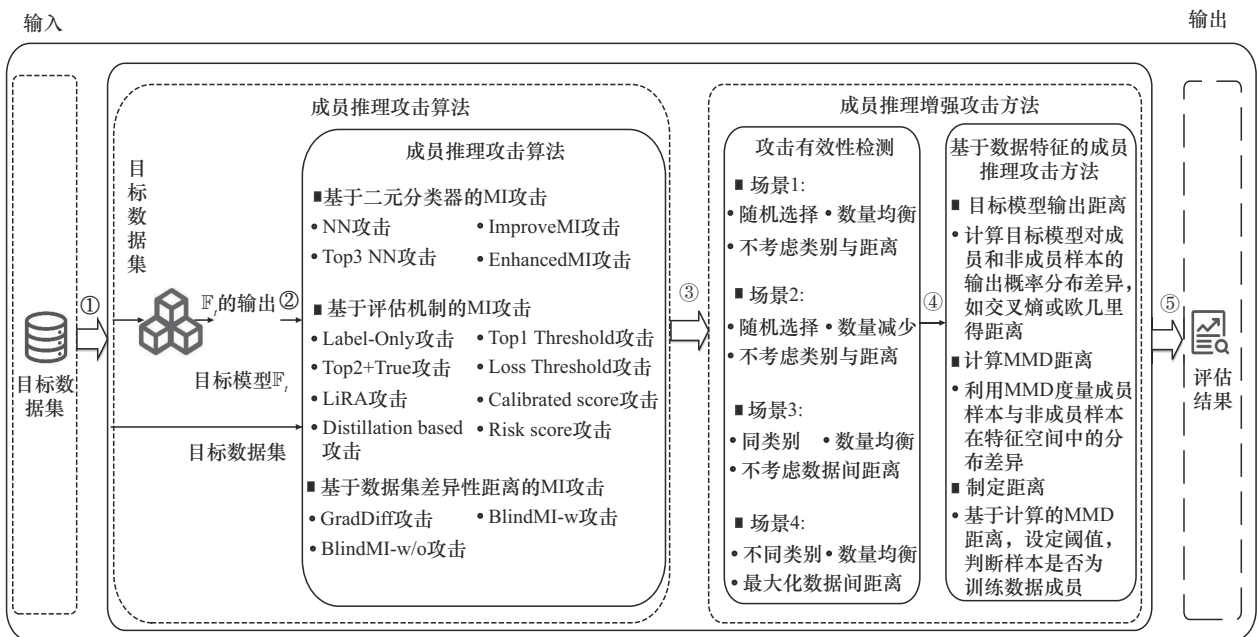


图1 基于数据特征的成员推理增强攻击方法

## 2.1 成员推理攻击有效性检测算法设计

### 1) 对数据样本进行分类

为了检测现有 MIA 的有效性, 本文首先将数据集中的样本分为 8 类: 正确分类的样本、错误分类的样本、成员样本、非成员样本、正确分类的成员样本、正确分类的非成员样本、错误分类的成员样本和错误分类的非成员样本。然后计算这 8 类样本分别对应的数量, 以及攻击的精确率 (precision)、召回率 (recall, 也称真阳率 (TPR, true positive rate))、F1-score、假阳率和成员增益。最后, 通过综合评估来判断 MIA 的有效性。

### 2) 计算数据样本间的距离

给定任意 2 个数据样本  $p$  和  $q$ , 首先将其输入目标模型获得输出概率。然后将所有输出概率映射到再生核希尔伯特空间<sup>[21]</sup>, 并计算数据样本间的最大均值差异性 (MMD, maximum mean discrepancy) 距离。

$$\text{Dis}[P_p, P_q] = \left| \frac{1}{n_p} \sum_{i=1}^{n_p} \phi(P_i) - \frac{1}{n_q} \sum_{j=1}^{n_q} \phi(P_j) \right|_{\mathcal{V}} \quad (1)$$

其中,  $\phi(P_i)$  和  $\phi(P_j)$  是再生核希尔伯特空间中的描述,  $\mathcal{V}$  是核空间的维度,  $P_i \in P_p$  且  $P_j \in P_q$ ,  $n_p(n_q)$  是样本  $p(q)$  的模型输出概率  $P_p(P_q)$  的大小。  $\phi(\cdot)$  是从概率空间到再生核希尔伯特空间的特征空间映射, 即  $k \mapsto \mathcal{V}$ , 高斯核函数最常用。

### 3) 成员推理攻击有效性检测算法

在对数据样本分类和计算距离后, 本文构造了 4 种测试场景, 并提出了一种成员推理攻击有效性检测算法。场景 1 (TS<sub>1</sub>) 的目标数据集  $D_{i1}$  分别在目标训练集和测试集 (或生成的非成员数据集) 中随机选择相等数量的成员和非成员 (如 10 000), 且不考虑数据集类别和数据样本间的距离。场景 2 (TS<sub>2</sub>) 的目标数据集  $D_{i2}$  分别在这 2 个数据集中随机选择相等数量且小于场景 1 数量的成员和非成员 (如 6 000), 且不考虑数据集类别和数据样本间的距离。场景 3 (TS<sub>3</sub>) 的目标数据集  $D_{i3}$  分别在这 2 个数据集中选择相等数量且为同一类别的成员和非成员 (如 6 000), 且考虑数据集类别, 但不考虑数据样本间的距离。场景 4 (TS<sub>4</sub>) 的目标数据集  $D_{i4}$  分别在这 2 个数据集中选择相等数量且不同类别的成员和非成员 (如 6 000), 并确保数据样本间距离最大。

算法 1 描述了成员推理攻击有效性检测算法。对于给定的 ML 模型、数据集和 MIA, 首先将数据集的数据样本分为 8 类, 计算每类的样本数量以及

数据样本间的距离。然后构造 4 个测试场景, 并检测已有 MIA 的有效性。算法 1 设置了一个假阳率阈值  $\alpha$  和一个成员增益阈值  $\varpi$ , 即若计算出的假阳率和成员增益满足阈值要求, 则认为该 MIA 有效。设置原理是现实中不同攻击者对假阳率的容忍度和成员增益的要求不同。因此, 本文将攻击者能容忍的最大的假阳率作为假阳率阈值  $\alpha$ , 将攻击者能接受的最小成员增益作为成员增益阈值  $\varpi$ 。

## 2.2 基于数据特征的成员推理增强攻击方法设计

利用成员推理攻击有效性检测算法检测现有 MIA 的有效性后, 为了降低攻击的假阳率并提高攻击的有效性, 本文根据攻击时使用的目标数据、数据特征和攻击原理, 提出了一种基于数据特征的成员推理增强攻击方法。数据特征主要包括目标模型对数据的输出概率、数据样本间的距离 (如 MMD 距离)。这些数据特征是数据自身的特征, 并没有涉及对数据处理后提取的特征。该方法的原理是基于模型对数据的输出概率以及输出概率之间的关联性 (如 MMD 距离), 来挖掘数据遭受攻击的风险, 进而选择合适的数据来增强攻击效果。

### 算法 1 成员推理攻击有效性检测

定义  $i = 1$ , 给定 ML 模型, 不同数据集  $D_{ii}$  和成员推理攻击 MIA <sub>$j$</sub> , 假阳率阈值  $\alpha$ , 成员增益阈值  $\varpi$

- 1) 对于第  $j$  个成员推理攻击 MIA <sub>$j$</sub>
- 2) 循环
- 3) 将数据集  $D_{ii}$  输入目标 ML 模型
- 4) 将数据集  $D_{ii}$  分为 8 类并计算每类的样本数量
- 5) 构造测试场景  $\text{TS} = \{\text{TS}_1, \text{TS}_2, \text{TS}_3, \text{TS}_4\}$
- 6) for  $\text{TS}_i \in \text{TS}$
- 7) 计算 precision、recall、F1-score、MA 和 FPR
- 8)  $i = i + 1$
- 9) 循环
- 10) until  $i = 4$ , 保存计算的 5 个评估指标数值
- 11) 循环
- 12) until  $j$  等于成员推理攻击的数量, 保存计算的评估指标的数值
- 13) 如果  $\text{MA} > \varpi$  且  $\text{FPR} < \alpha$
- 14) 返回攻击方法有效
- 15) 否则
- 16) 返回攻击方法无效
- 17) end for

**算法2** 基于数据特征的成员推理增强攻击

定义  $i = 1$ , 给定 ML 模型, 不同的目标数据集  $D_{ii}$  和成员推理攻击  $MIA_j$ , 最小假阳率  $\delta$ , 初始化可测试的目标数据集  $D_{test}$

- 1) 对于不同的目标数据集  $D_{ii}$
- 2) 循环
- 3) 将数据集  $D_{ii}$  输入目标 ML 模型
- 4) 根据 ML 模型的输出将数据集  $D_{ii}$  分为两类
- 5) 用式 (3) 计算数据样本间的 MMD 距离  $MMD_{hl}$
- 6) 对于第  $j$  个成员推理攻击  $MIA_j$
- 7) 从算法 1 获得假阳率为  $\delta$  对应的 MMD 距离, 记作距离阈值  $\pi$
- 8) 如果  $MMD_{hl}$  大于  $\pi$ , 即  $MMD_{hl} > \pi$  成立
- 9) 则算法收敛, 将此时的目标数据集存入可测试的目标数据集, 即  $D_{test} \leftarrow \{D_{ii}\}$
- 10) 否则, 算法未收敛, 利用数据生成技术生成一些非成员数据
- 11) 将生成的非成员数据与“输出概率低的样本”进行替换, 重新组成一个新的测试数据集  $D_{ii}^n$
- 12) 用式(3)计算  $MMD_{hm}$
- 13) 如果  $MMD_{hm}$  大于  $\pi$ , 即  $MMD_{hm} > \pi$  成立
- 14) 则算法收敛,  $D_{test} \leftarrow \{D_{ii}^n\}$
- 15) 返回可测试的目标数据集  $D_{test}$

算法2描述了基于数据特征的成员推理增强攻击方法, 对于一个目标数据集, 首先根据目标模型对目标数据的输出概率将数据分为两类——“输出概率高的样本”和“输出概率低的样本”。接着用式(1)计算这两类数据样本间的MMD距离  $MMD_{hl}$ , 并根据攻击者能容忍的最小假阳率  $\delta$ , 利用算法1计算该

最小假阳率下数据样本间的MMD距离, 将其作为距离阈值  $\pi$ , 若  $MMD_{hl}$  大于  $\pi$ , 则说明假阳率小于  $\delta$ , 即该目标数据集可用。否则, 利用数据生成技术<sup>[7]</sup>生成一些非成员数据, 并将该非成员数据和“输出概率低的样本”进行替换, 然后再计算“输出概率高的样本”和生成的非成员数据 MMD 距离  $MMD_{hm}$ , 并判断  $MMD_{hm}$  与距离阈值  $\pi$  的相对大小, 该操作一直进行到  $MMD_{hm}$  大于  $\pi$ 。由于目标数据集中数据样本的成员关系无法事先预知, 本文首先根据模型输出概率将数据样本分为上述两类, 同时已有研究<sup>[3,5]</sup>表明目标模型对成员样本的表现优于非成员样本, 因此本文认为“输出概率高的样本”更有可能是成员, “输出概率低的样本”更有可能是非成员。在对数据样本进行替换时, 首先将“输出概率高的样本”的索引固定住, 只将生成样本与“输出概率低的样本”进行替换, 并确保替换后数据样本与“输出概率高的样本”间的距离大于距离阈值  $\pi$ 。最终实现攻击者可接受的最小假阳率下的攻击效果。

**3 实验评估**

本节首先介绍了实验设置、模型架构、评估指标以及数据集, 然后评估了本文所提鲁棒性检测方法和攻击增强方法的有效性。

**3.1 实验设置**

本文实验的操作系统是 Ubuntu 22 64 位, CPU 是 AMD EPYC 7542 32-Core Processor, GPU 是 NVIDIA GeForce RTX 3090, CUDA 版本是 12.0, Python 版本是 3.10, Pytorch 版本是 2.3.0。

**3.2 模型架构和评估指标**

如表 3 所示, 参考 BlindMI 攻击<sup>[7]</sup>本文实验采用了多层感知器 (MLP, multilayer perceptron) 模型,

**表 3** 目标模型和影子模型架构以及超参数设置

模型架构	层数	目标模型		影子模型	
		最大的 Epoch	学习率	最大的 Epoch	学习率
MLP	3~7	$e_m$	$5 \times 10^{-5}$	$[0.3\sim 2]e_m$	$1 \times 10^{-4}$ 或 $1 \times 10^{-5}$
StandDNN	2	$e_m$	$5 \times 10^{-5}$	$0.5e_m$	$1 \times 10^{-4}$
VGG16	16	$e_p + e_m$	$5 \times 10^{-5}$	$e_p + 0.6e_m$	$5 \times 10^{-5}$
VGG19	19	$e_p + e_m$	$5 \times 10^{-5}$	$e_p + 1.5e_m$	$5 \times 10^{-5}$
ResNet50	50	$e_p + e_m$	$5 \times 10^{-5}$	$e_p + 0.2e_m$	$5 \times 10^{-5}$
ResNet101	101	$e_p + e_m$	$5 \times 10^{-5}$	$e_p + 0.3e_m$	$1 \times 10^{-4}$
DenseNet121	121	$e_p + e_m$	$5 \times 10^{-5}$	$e_p + e_m$	$1 \times 10^{-4}$

最多有 7 个密集层 (8 192、4 096、2 048、1 024、512、256 和 128) 和一个 Softmax 层, 其中,  $e_m$  表示目标模型在表 4 的每个数据集上的最大 Epoch,  $e_p$  表示在 ImageNet 数据集上预训练权重的 Epoch。标准的卷积神经网络 (CNN, convolutional neural network) 架构和超参数与已有工作<sup>[3]</sup>相同, 还利用了 3 种流行的深度神经网络 (DNN, deep neural network) 模型。给定一个训练集, 实验随机选择一个目标模型, 并用指定超参数训练该目标模型。黑盒攻击随机选择一个影子模型架构, 灰盒和白盒攻击需选择与给定目标模型相同的影子模型。

### 3.3 数据集

如表 4 所示, 参考已有工作<sup>[7,18]</sup>, 本文实验采用 MIA 领域 4 个广泛使用数据集: CIFAR100、CIFAR10、CH\_MNIST 和 ImageNet, 且测试集中成员和非成员数量相等。

### 3.4 攻击有效性检测算法的有效性

本文在 4 个测试场景上测试各个攻击的效果, 当测试的假阳率小于  $\alpha$  且成员增益大于  $\varpi$  时, 将其视为攻击有效, 反之亦然。本文实验设置  $\alpha$  和  $\varpi$  的取值范围分别是 {5%, 10%, 15%, 20%} 和 {20%, 30%, 40%, 50%}。从图 2~图 9 可知, 尽管已有

表 4 实验中使用的数据集

数据集	场景	类别	分辨率	Epoch 大小	训练数据集	参考数据集	测试集成员	测试集非成员
CIFAR10	图像分类	10	32×32	(pre-trained + 30)或 150	25 000	25 000	12 500	12 500
CIFAR100	图像分类	100	32×32	(pre-trained + 30)或 150	25 000	25 000	12 500	12 500
CH_MNIST	图像分类	8	64×64	(pre-trained + 15)或 150	2 250	2 250	500	500
ImageNet	图像分类	200	64×64	(pre-trained + 30)或 150	10 000	10 000	20 000	20 000

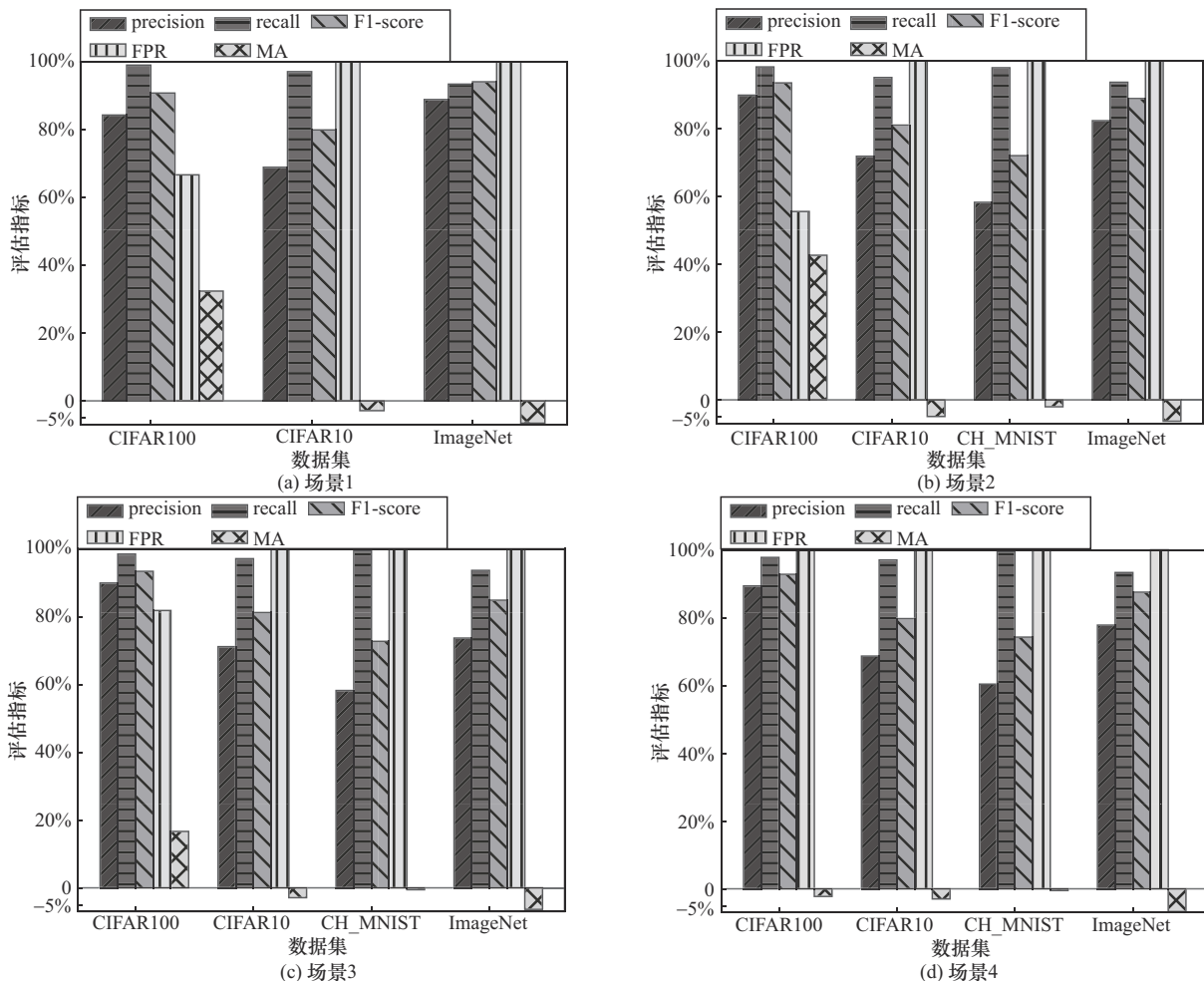


图 2 BlindMI-w/o 攻击在 4 种测试场景中的攻击效果

MIA在不同测试场景下的F1-score(如87.49%)和召回率(如99.17%)较高,但是其假阳率(如100%)却远远超出了假阳率阈值 $\alpha$ ,进一步限制了其实际应用。例如,图2(a)~图2(d)分别展示了BlindMI-w/o攻击在场景1、场景2、场景3和场景4中的攻击效果,可以看出虽然BlindMI-w/o攻击的召回率和F1-score较高(如99.04%),但其假阳率最高可达100%,而攻击者能接受的假阳率通常为10%或者更低。因此,BlindMI-w/o攻击的实际应用不是很理想。图3表明,尽管BlindMI-w攻击在CIFAR10和CH\_MNIST数据集上的召回率和F1-score较高(如92.77%),但其仍存在较高的假阳率(如75%)。GradDiff攻击也基于数据集差异来实施攻击,表现出了与BlindMI-w攻击相似的攻击效果。

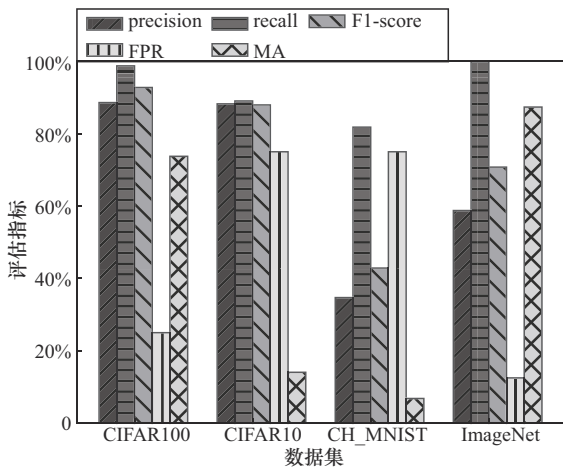


图3 BlindMI-w攻击在测试场景1下的攻击效果

由图4可知,当测试数据集的数据量发生变化时,Label-Only攻击仍表现出较高的假阳率。例

如,当图4(a)中测试数据为10 000时,Label-Only攻击在CH\_MNIST上的假阳率达到了81.04%。当图4(b)中测试数据为6 000时,Label-Only攻击在CIAFR10和CIFAR100数据集上的假阳率都达到了100%,而实际中这么高的假阳率是不可用的。图5表明,尽管Loss Threshold攻击的召回率(如99.07%)和F1-score(如87.35%)较高,但是其精确率(如51.46%)和成员增益(如9.95%)较低。经分析发现Loss Threshold攻击的假阳率较高,如在测试场景1和场景2上,Loss Threshold攻击在ImageNet数据集上的假阳率分别达到了94.77%和100%,在实际中是不可接受的。观察图6可知,NN攻击在ImageNet数据集上召回率达到了99.42%,但其精确率和成员增益较低(分别为50.05%和-3.58%),究其原因发现其假阳率达到了100%,极大地损害了攻击的有效性。同时可发现当测试数据集的大小发生变化时,对其攻击有效性的影响较小。EnhancedMI和ImproveMI攻击的攻击效果与NN攻击基本一致。

图7表明,Top1 Threshold攻击的F1-score最高可达79.59%,但其成员增益却只有18.63%,而假阳率高达93.33%,制约其实际可行性。由图8可知,只关注Top2+True攻击的召回率和F1-score时(如98.36%和70.79%),可认为其是有效的,但其精确率和成员增益分别为55.30%和18.84%。因此实际中Top2+True攻击并非有效。图9表明Top3 NN攻击的F1-score和召回率较高(91.12%),而其成员增益较低(26.24%),原因在于不同测试场景下该攻击的假阳率较高(如100%),超出了承受范

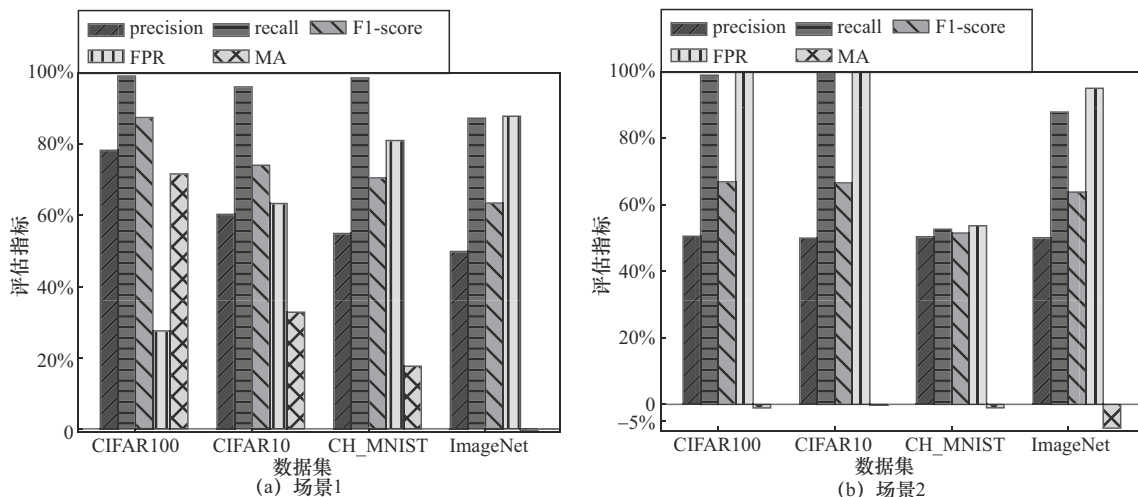


图4 Label-Only攻击在测试场景1和场景2中的攻击效果

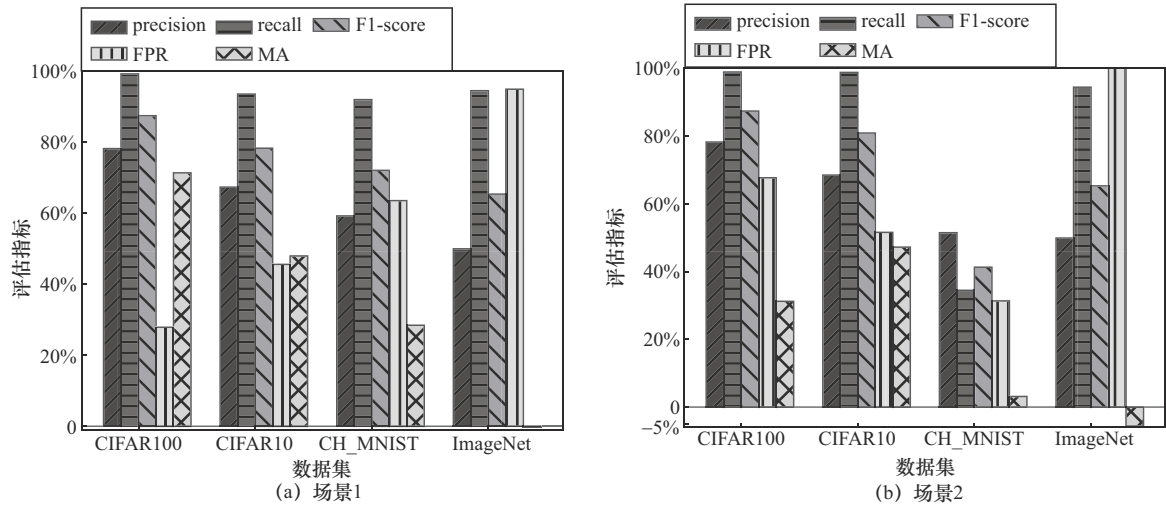


图5 Loss Threshold攻击在测试场景1和场景2中的攻击效果

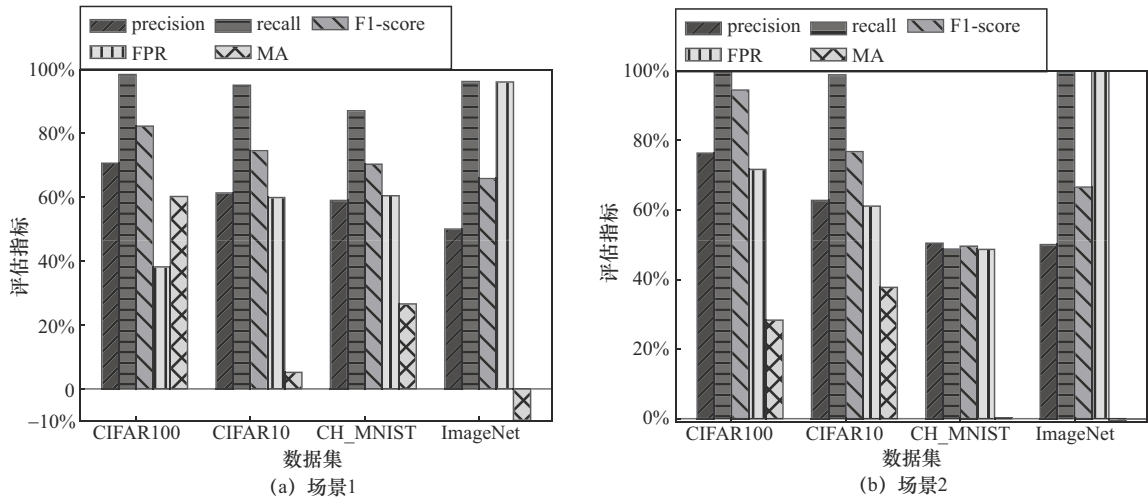


图6 NN攻击在测试场景1和场景2中的攻击效果

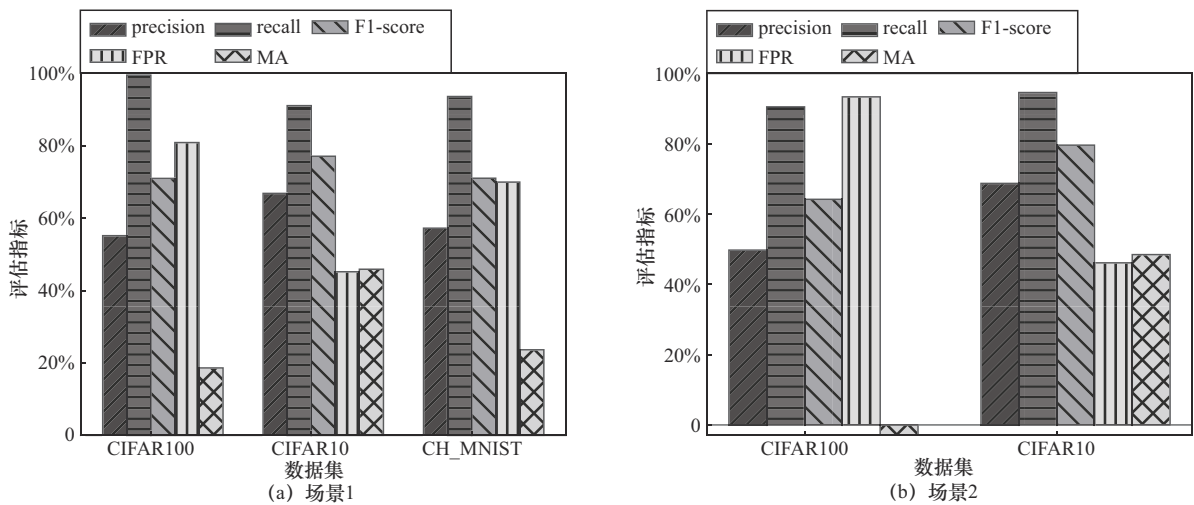


图7 Top1 Threshold攻击在测试场景1和场景2中的攻击效果

围。同时，观察表5可知，不同攻击的曲线下面积 (AUC, area under the curve) 和TPR@0.1%FPR这两个

指标的取值均较低，也从侧面揭示了已有工作有较高的假阳率。

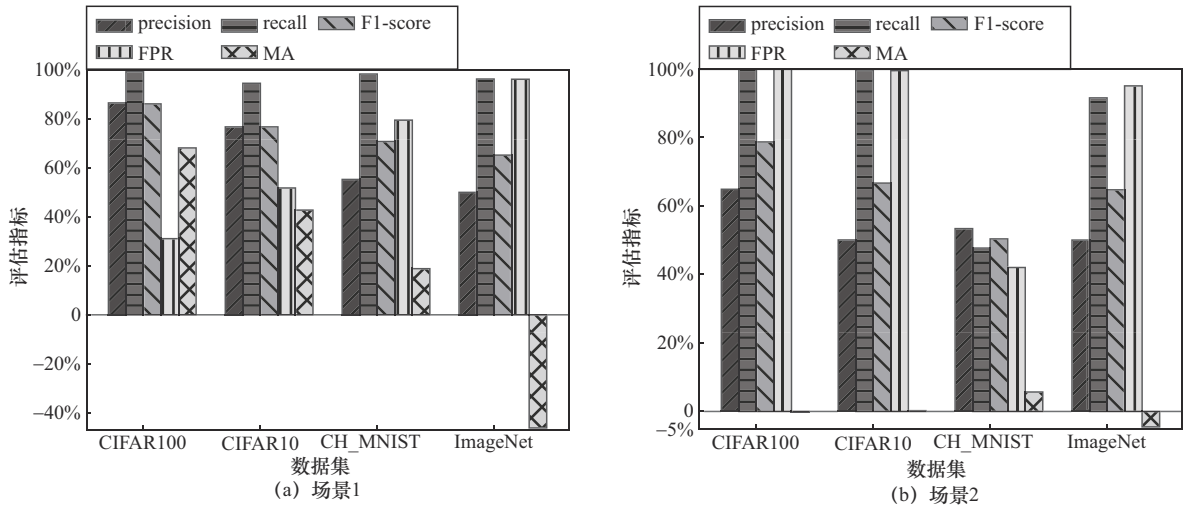


图 8 Top2+True 攻击在测试场景 1 和场景 2 中的攻击效果

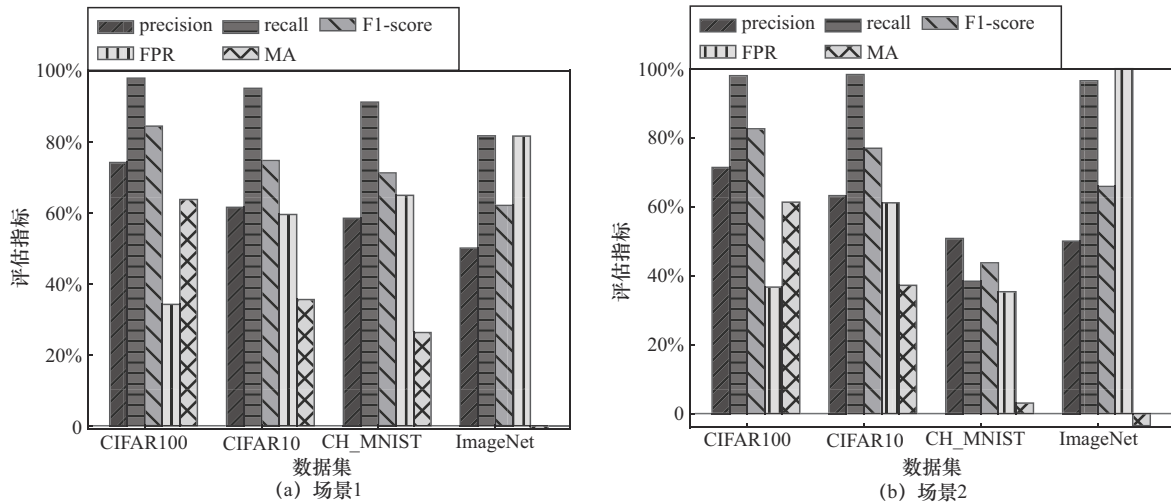


图 9 Top3 NN 攻击在测试场景 1 和场景 2 中的攻击效果

表 5 不同攻击方法 AUC 和 TPR@0.1%FPR 的测试结果

攻击方法	AUC	TPR@0.1%FPR
NN	0.53	0.12%
Top3 NN	0.52	0
Label-Only	0.55	0
Top1 Threshold	0.828	0
BlindMI-w	0.501	0
Top2+True	0.597	0

### 3.5 基于数据特征的成员推理增强攻击的有效性

本节主要在 4 个数据集上分别对比了利用本文所提的增强攻击方法处理前 (O, original) 和处理后 (R, reduction), 15 种 MIA 的攻击效果。图 10 表明, 本文方法在降低 BlindMI-w/o 和 BlindMI-w 攻击假阳率的同时, 将这 2 个攻击的成员增益分别提升了

67.63% 和 36.4%。由此可见, 本文方法对 BlindMI-w/o 攻击的性能提升力度更大, 主要原因是 BlindMI-w/o 攻击直接利用测试数据集作为非成员数据集来捕捉成员关系, 而 BlindMI-w 攻击则利用生成的非成员数据集。而生成的非成员数据与替换数据都是生成数据 (可能采用相同的数据生成技术), 导致替换后数据间的距离相较于测试数据与生成数据的有所减小, 从而限制了 BlindMI-w 攻击性能提升。同时, GradDiff 与 BlindMI-w 攻击的性能提升幅度基本一致。

图 11 表明在 CIFAR100 数据集上, 相较于 Loss Threshold 攻击, 本文方法对 Label-Only 攻击的性能增强力度更大。例如, 将 Label-Only 攻击精确率和成员增益分别提升了 10.38% 和 55.27%。原因是 Label-Only 攻击利用预测标签与实际标签的一致性来

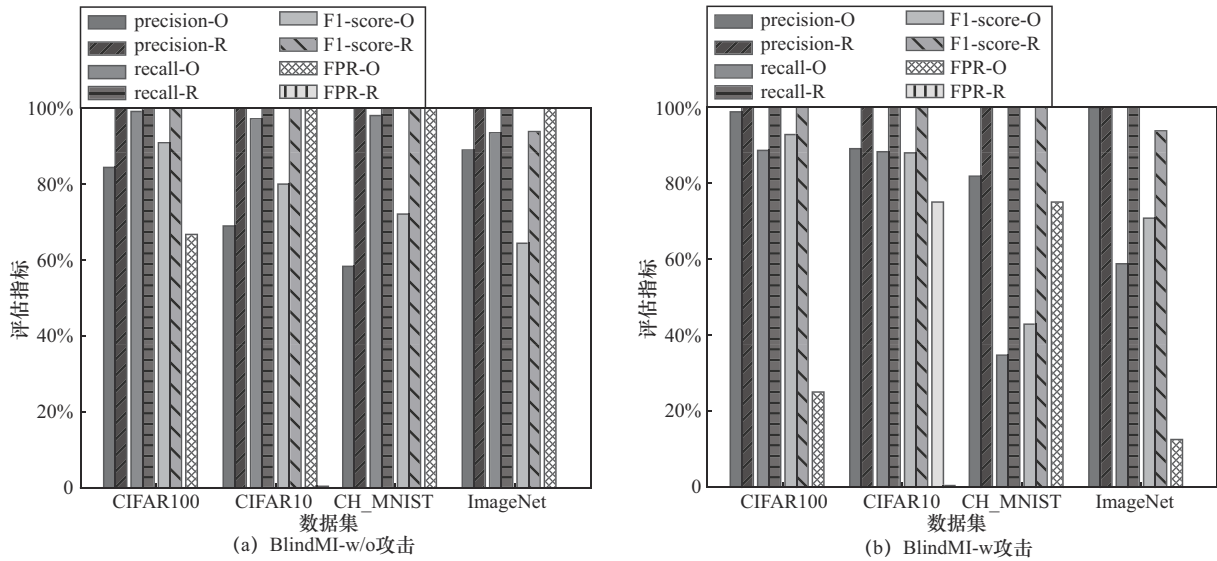


图10 BlindMI-w/o和BlindMI-w攻击处理前和处理后攻击效果对比

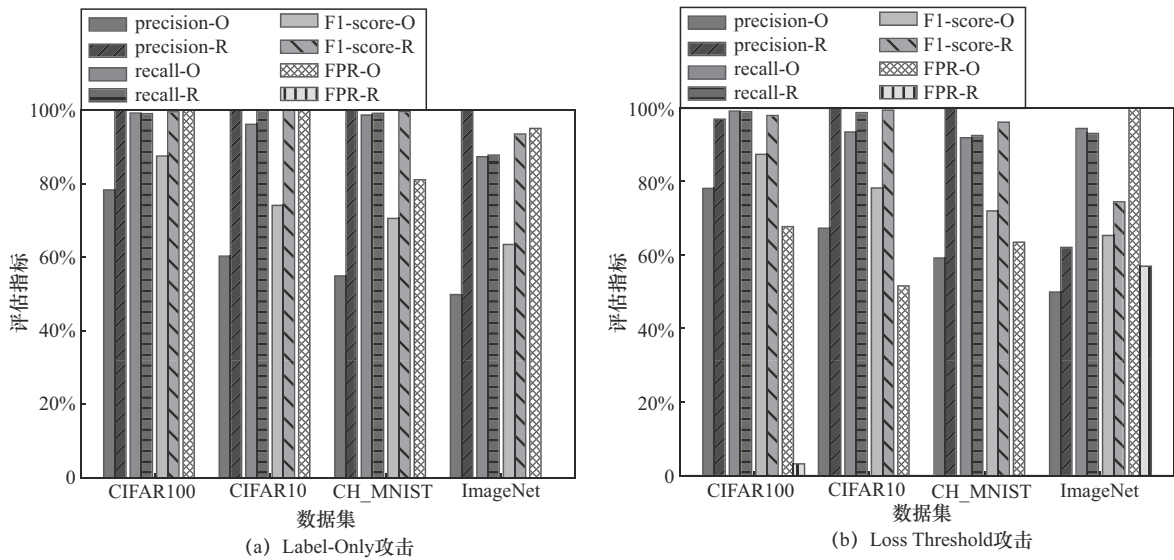


图11 Label-Only和Loss Threshold攻击处理前和处理后攻击效果对比

识别成员，Loss Threshold攻击将目标样本的交叉熵损失与成员样本的平均交叉熵损失的相对大小作为成员依据。在用生成的非成员数据对测试数据中“输出概率低的样本”进行替换时，由于替换后数据标签保持不变，且替换后数据距离增大了，降低了样本间的区分难度，极大增强了Label-Only攻击的性能。而替换的样本虽确保了数据间的差异性，但无法获取成员样本的交叉熵损失，因而对Loss Threshold攻击的性能提升相对较小。

如图12所示，相较于Top1 Threshold攻击，本文方法可将NN攻击的假阳率降低幅度较大，而对这2个攻击的成员增益增强幅度大致相当（如

46.91%和41.47%）。主要原因是NN攻击是基于模型的全部输出概率来实施攻击，替换后的样本保证了测试数据间的距离最大，增大了数据样本间的差异，促使该攻击更易捕捉成员关系，进而提升了攻击性能。而替换的样本很难精准刻画Top1非成员阈值，进而对Top1 Threshold攻击的假阳率降低和性能增强幅度相对NN攻击低。由于EnhancedMI和ImproveMI攻击与NN攻击的原理一致，因此这2个攻击增强效果与NN攻击的基本一致。

图13表明，相对于Top3 NN攻击，本文方法对Top2+True攻击的成员增益增强幅度稍大。经分析发现本文方法可对Top2+True攻击的F1-score增

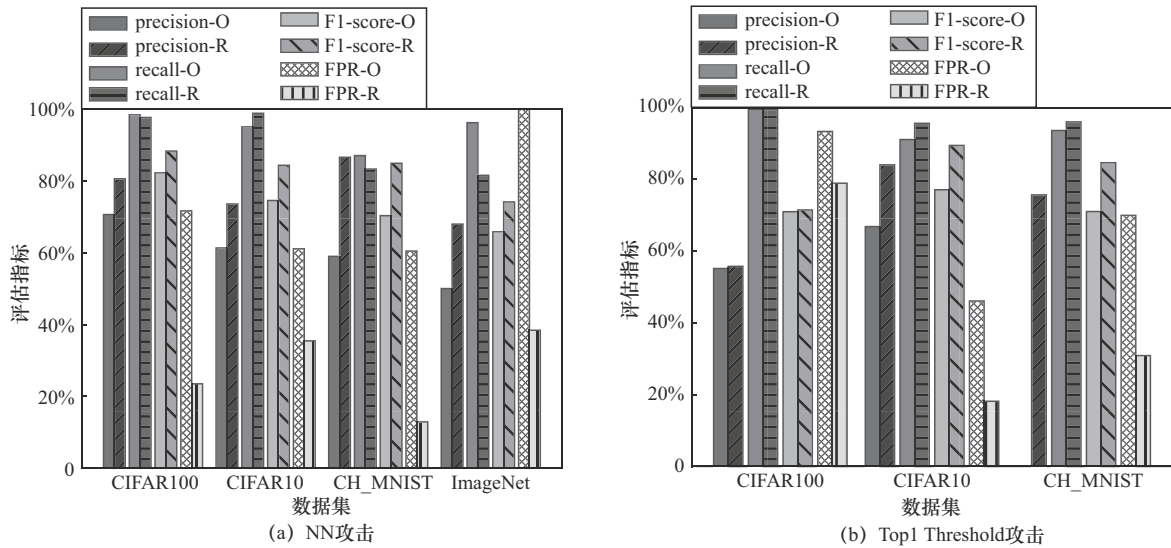


图 12 NN和Top1 Threshold攻击处理前和处理后攻击效果对比

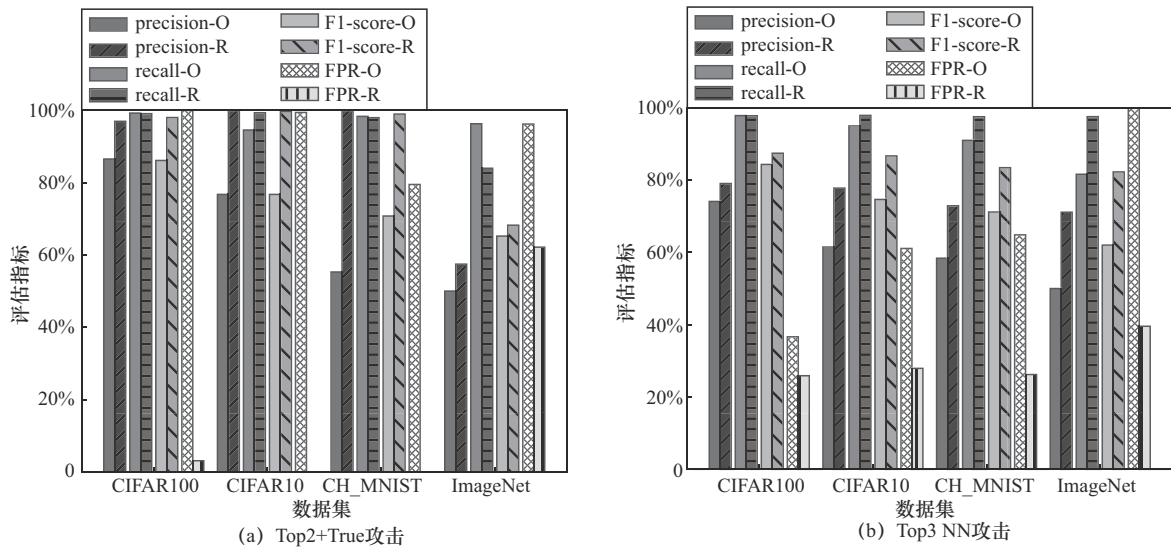


图 13 Top2+True和Top3 NN攻击处理前和处理后攻击效果对比

强幅度（如 28.2%）大于对 Top3 NN 攻击的增强幅度（如 20.25%）。主要原因是替换的样本会改变测试样本间的距离，对这 2 个攻击均有影响。又由于 Top3 NN 攻击仅利用 Top3 输出概率，而 Top2+True 攻击不仅采用了 Top2 输出概率还联合了真实标签，因此本文方法对 Top2+True 攻击的假阳率降低幅度（如 99.44%）相较 Top3 NN 攻击较大（如 60.42%）。此外，Top2 相较于 Top3 输出概率限制较宽松，从而进一步增强了 Top2+True 攻击的性能。

如图 14 所示，相较于 LiRA 攻击，Distillation based 攻击精确率增大幅度提升了 91.31%。主要原因是 Distillation based 攻击同时依赖模型和数据的攻击阈值来进行成员关系识别，而本文方法不仅利

用数据生成技术生成测试数据来捕捉数据特征。而且对“输出概率低的样本”进行替换，并确保替换后模型输出的数据样本间距离较大，其进一步捕捉了模型特征，使构建的攻击阈值更精确，进而对 Distillation based 攻击的性能提升较大。而 LiRA 攻击联合了样本的难度校准分数和高斯似然估计来识别成员关系。尽管替换的样本可以在一定程度上提升该攻击性能，但由于原始训练集无法预知，使对 LiRA 攻击性能提升相对较低。

图 15 表明，相较于 Calibrated Score 攻击，Risk score 攻击的 F1-score 增大幅度提升了 58.28%。主要原因是 Risk score 攻击联合样本的隐私风险分数和预测熵来实施攻击。本文方法仅对“输出概率

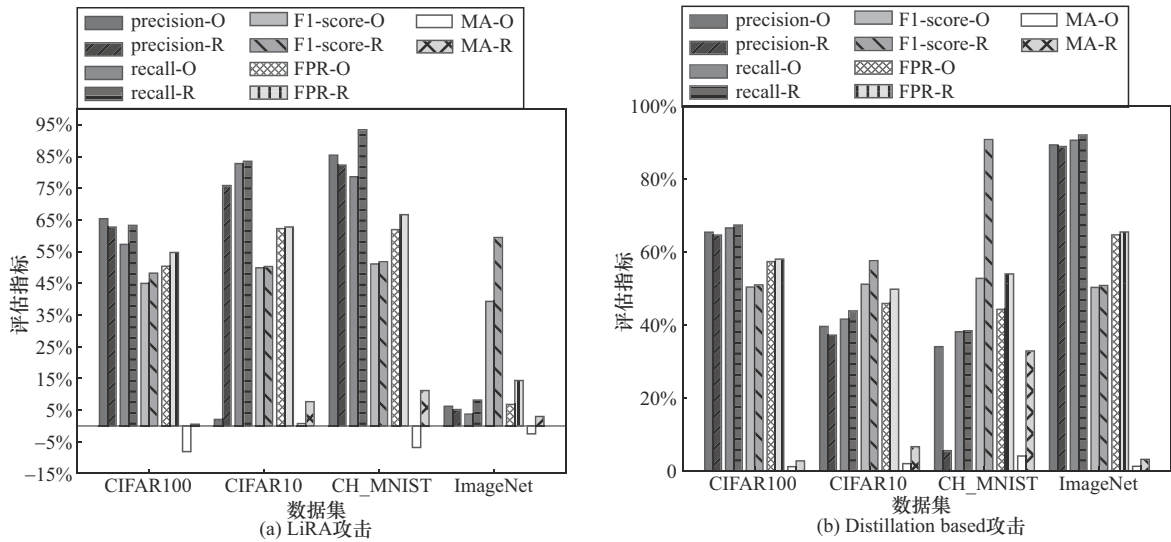


图 14 LiRA 和 Distillation based 攻击处理前和处理后攻击效果对比

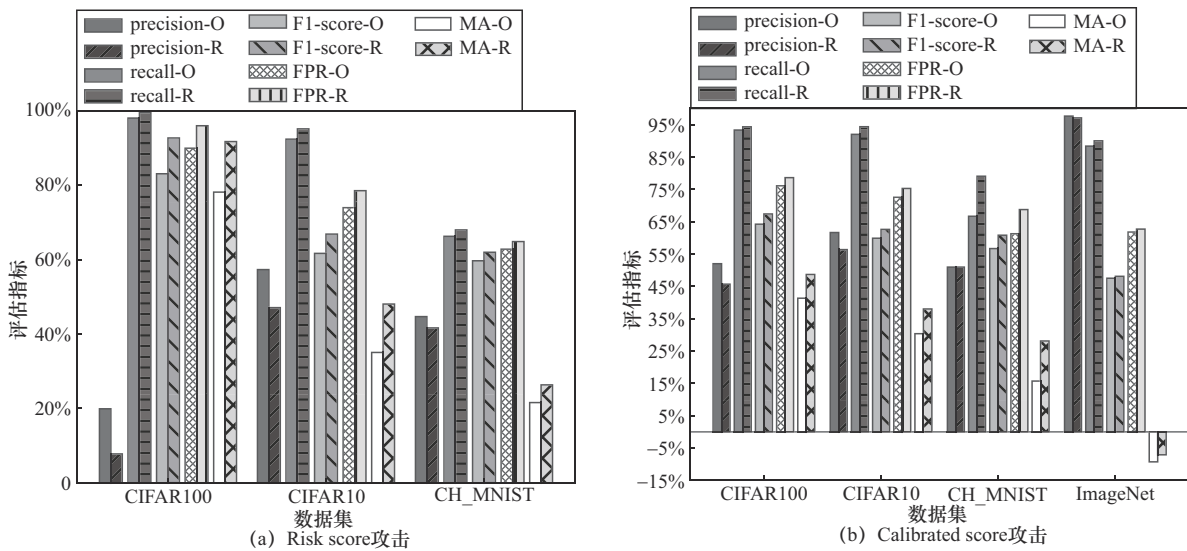


图 15 Risk score 和 Calibrated score 攻击处理前和处理后攻击效果对比

低的样本”进行替换，并确保替换后数据样本间的距离较大，增大了测试数据集中隐私风险高的目标样本与低隐私风险样本间的差异，进而增强其攻击性能。而替换后的数据可辅助 Calibrated score 攻击计算更精确的难度校准分数，又由于训练集未知，因此对 Calibrated score 攻击的性能提升相对较小。

### 3.6 消融实验

由图 16 可知，当数据集需要替换的数据增多时，完成一次攻击所需的时间也有所增加。由于本文仅根据目标模型对数据的输出概率来计算 MMD 距离，并样本进行替换和攻击。因此，数据替换对攻击整体的执行时间影响不大。图 17 表明，随着数据替换比例增加，攻击效果逐渐增强。主要原因

是替换的数据样本越多，数据样本间的距离越大，攻击越容易捕捉和识别成员和非成员的特征。

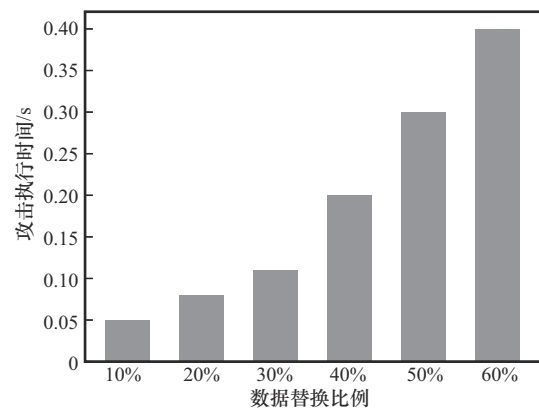


图 16 数据替换比例与攻击执行时间的关系

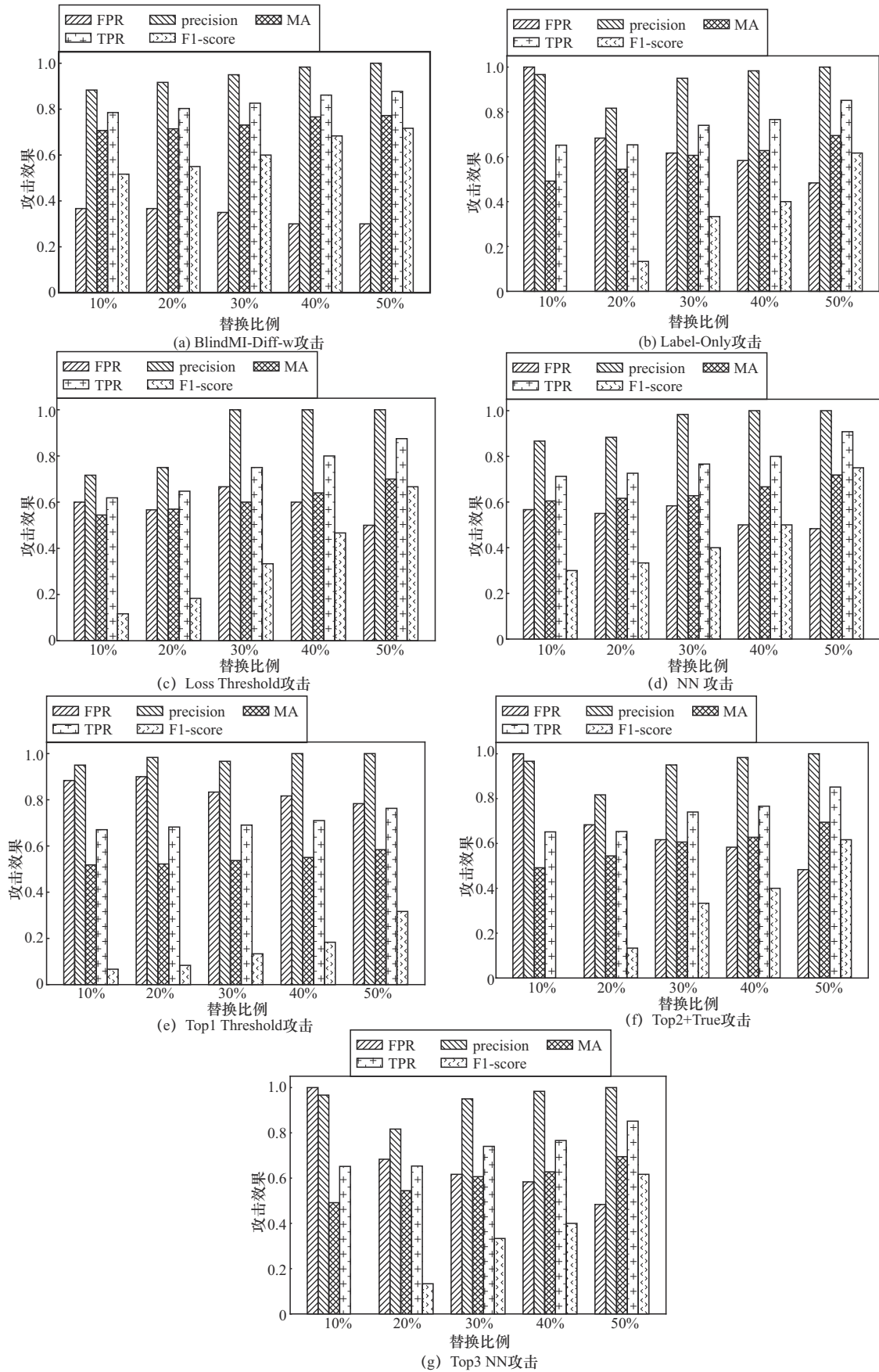


图17 对数据集进行不同比例的替换后对攻击效果的影响

### 4 方法讨论

上述研究表明，已有 15 种 MIA 仅仅在场景 1 和场景 2 下测试攻击效果，并未考虑场景 3 和场景 4 的情况，本文在构造的 4 个测试场景上发现现有 15 种 MIA 的假阳率最高可达 100%。此外，实验还

观察到相比于场景 1 和场景 2，本文所提的有效性检测算法在检测场景 3 上的假阳率更高。例如，在 CIFAR100 数据集的场景 1 和场景 3 上，BlindMI-w/o 攻击的假阳率分别为 66.67% 和 81.82%。图 18~图 21 展示了 15 种 MIA 在原始场景（没有任何改进）和改进场景（本文方法）下的攻击效果。

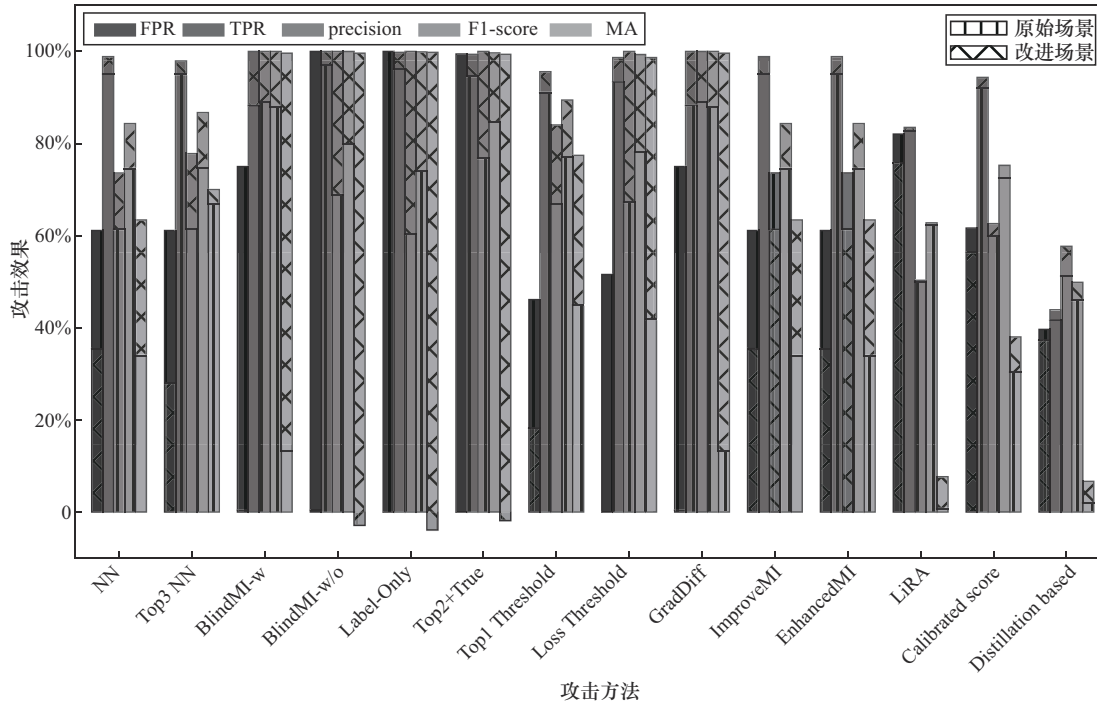


图 18 CIFAR100 上 15 种 MIA 在原始场景和改进场景下的攻击效果

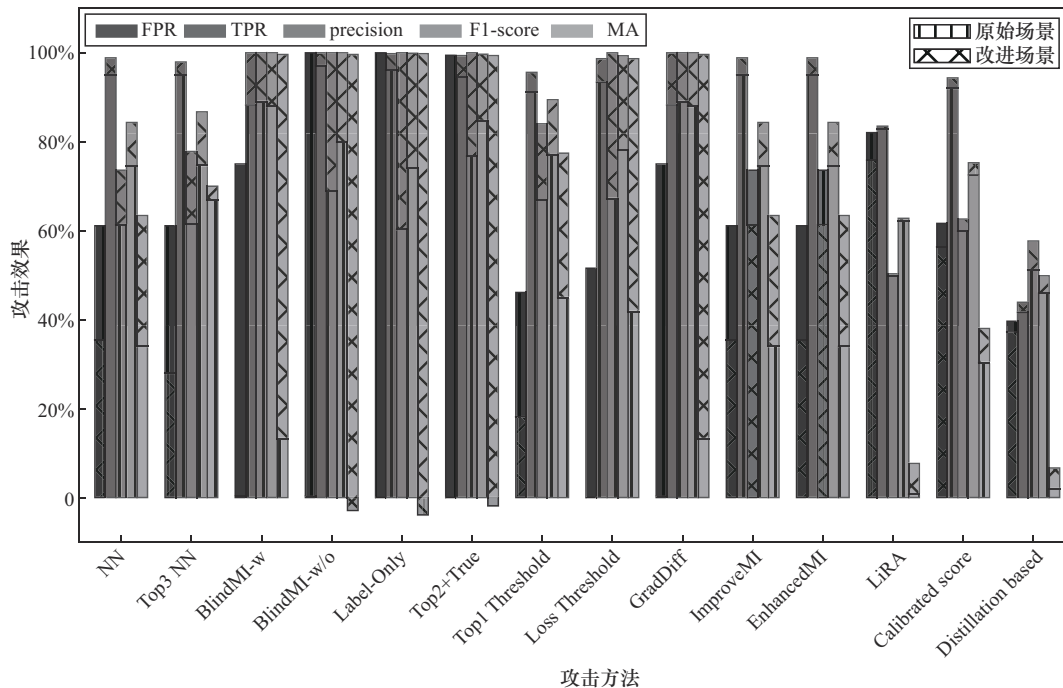


图 19 CIFAR10 上 15 种 MIA 在原始场景和改进场景下的攻击效果

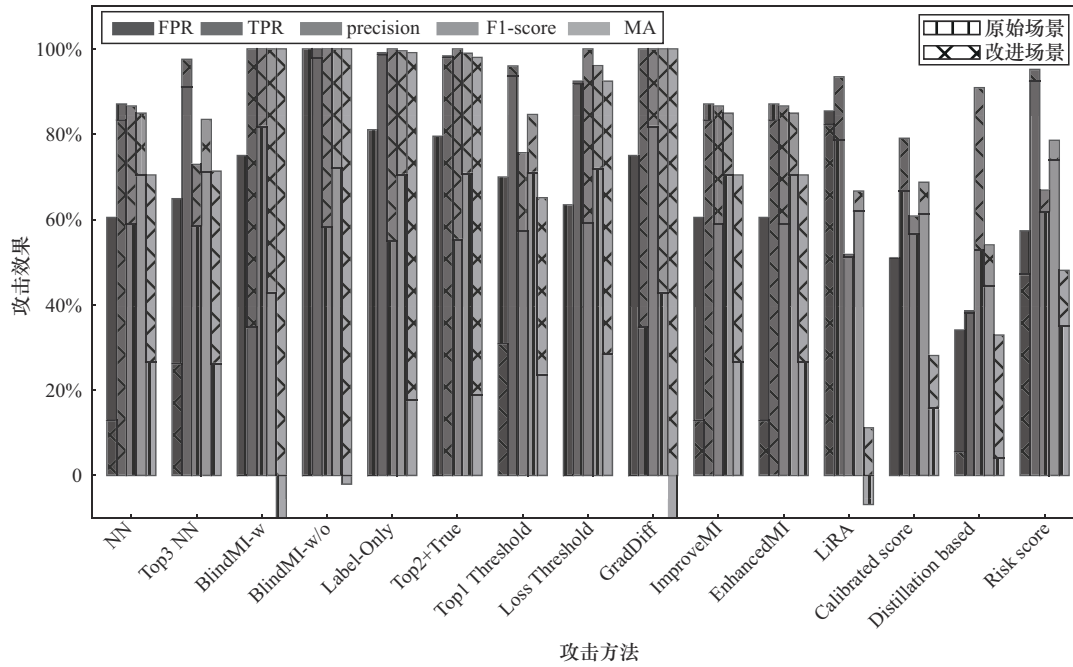


图 20 CH\_MNIST 上 15 种 MIA 在原始场景和改进场景下的攻击效果

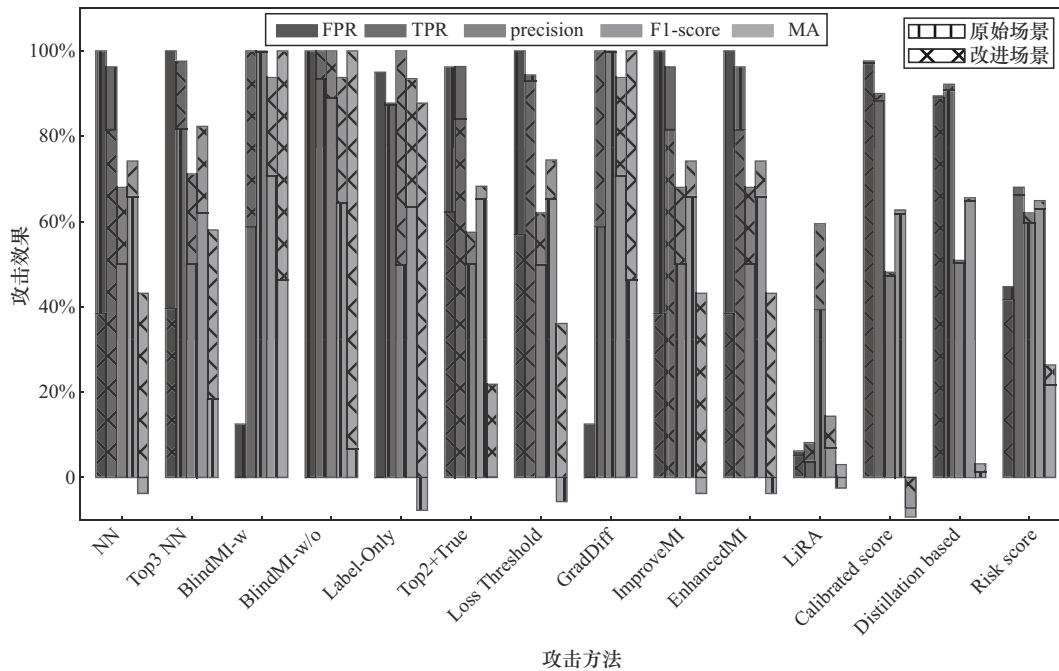


图 21 ImageNet 上 15 种 MIA 在原始场景和改进场景下的攻击效果

由图 18~图 21 可知, 相比于原始场景, 在利用本文方法改进后的场景上, 15 种 MIA 在 4 个数据集上的攻击效果有显著提升。例如, 在 CIFAR100 数据集上, 可将 BlindMI-w 攻击的精确率、召回率、F1-score、成员增益分别从 84.31%、99.04%、90.80% 和 32.37% 提高到 100%、100%、100% 和 100%, 将假阳率从 66.67% 降低到 0。主要原因是

本文主要基于数据自身的特征 (如目标模型对数据的输出概率以及数据样本间的距离) 来对目标数据集集中的样本进行筛选, 并利用满足假阳率和成员增益条件的目标样本来实施攻击。依据的原理是模型对训练数据和测试数据的表现不同, 模型更容易记住训练数据, 进而在训练数据上表现更好, 在测试数据上表现较差。已有工作<sup>[17,19]</sup>表明, 目标数据集

中存在和成员样本相类似的非成员样本，其大大增加了攻击的假阳率。同时，MIA 通常仅能获取目标模型输出结果，而无法访问原始数据标签。MMD 距离的无监督特性使其特别适用于此类攻击。此外，尽管其他距离度量指标（如 Kullback-Leibler 散度），也可以用来测量分布差异，但由于存在计算不稳定或对噪声敏感等问题，使其在高维空间和复杂分布场景中表现不佳。因此，本文利用 MMD 距离可准确评估不同样本间的相似性，而样本间的 MMD 距离越大，降低了样本区分难度，从而增强攻击效果。

本文提出的成员推理攻击有效性检测算法虽可检测已有 MIA 的有效性。但由于不同攻击者对假阳率的可接受能力不同，因此该算法设置的假阳率阈值  $\alpha$  和成员增益阈值  $\varpi$  很难统一，对于不同的攻击者需要进行多次测试。此外，基于数据特征的成员推理增强攻击方法主要利用数据特征筛选目标样本，这些数据特征包括目标模型对数据的输出概率以及数据样本间的距离，其并没有涉及对数据处理后提取的特征。因此，该方法仅能提升基于模型输出概率的攻击性能，而对于基于其他原理的攻击的增强效果仍不清楚。同时，当多个攻击者（如  $l$  个）利用该方法来增强攻击性能时，由于不同场景下不同攻击者对假阳率要求不同。因此，当其达到假阳率要求时需测试的次数分别为  $\{\text{Num}_1, \text{Num}_2, \dots, \text{Num}_l\}$ ，那么完成一次增强攻击所需的测试次数为这  $l$  个次数之和。

## 5 结束语

针对已有工作未充分检测现有成员推理攻击方法的有效性和实用性，导致错误推测误导攻击者的判断，本文首先考虑数据特征设计了一种成员推理攻击有效性检测算法，并检测出现有攻击方法的假阳率高达 100%。然后，结合攻击的基本原理，提出了一种基于数据特征的成员推理增强攻击方法，在降低攻击假阳率的同时提高攻击效率。本文方法在多个基准数据集上，可将 15 个现有 MIA 方法的假阳率降低到 0，同时将攻击精确率和成员增益都提高到 100%，实现了攻击效果的增强。

## 参考文献：

[1] 郭虎升, 孙妮, 王嘉豪, 等. 基于自适应深度集成网络的概念漂移收敛

方法[J]. 计算机研究与发展, 2024, 61(1): 172-183.

GUO H S, SUN N, WANG J H, et al. Concept drift convergence method based on adaptive deep ensemble networks[J]. Journal of Computer Research and Development, 2024, 61(1): 172-183.

[2] 廖鹏, 方滨兴, 刘潮歌, 等. NFT 仿冒欺诈的测量与检测技术[J]. 计算机学报, 2024, 47(5): 1065-1081.

LIAO P, FANG B X, LIU C G, et al. Measurement and detection techniques for NFT counterfeiting fraud[J]. Chinese Journal of Computers, 2024, 47(5): 1065-1081.

[3] SHOKRI R, STRONATI M, SONG C Z, et al. Membership inference attacks against machine learning models[C]//Proceedings of the 2017 IEEE Symposium on Security and Privacy (SP). Piscataway: IEEE Press, 2017: 3-18.

[4] LIU H, WU Y H, YU Z Y, et al. Please tell me more: privacy impact of explainability through the lens of membership inference attack[C]//Proceedings of the 2024 IEEE Symposium on Security and Privacy (SP). Piscataway: IEEE Press, 2024: 4791-4809.

[5] SALEM A, ZHANG Y, HUMBERT M, et al. ML-leaks: model and data independent membership inference attacks and defenses on machine learning models[C]//Proceedings 2019 Network and Distributed System Security Symposium. Internet Society, 2019: 1-15.

[6] YEOM S, GIACOMELLI I, FREDRIKSON M, et al. Privacy risk in machine learning: analyzing the connection to overfitting[C]//Proceedings of the 2018 IEEE 31st Computer Security Foundations Symposium (CSF). Piscataway: IEEE Press, 2018: 268-282.

[7] HUI B, YANG Y C, YUAN H L, et al. Practical blind membership inference attack via differential comparisons[C]//Proceedings 2021 Network and Distributed System Security Symposium. Internet Society, 2021: 1-17.

[8] NASR M, SHOKRI R, HOUMANSADR A. Machine learning with membership privacy using adversarial regularization[C]//Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security. New York: ACM Press, 2018: 634-646.

[9] ABADI M, CHU A, GOODFELLOW I, et al. Deep learning with differential privacy[C]//Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security. New York: ACM Press, 2016: 308-318.

[10] JIA J Y, SALEM A, BACKES M, et al. MemGuard: defending against black-box membership inference attacks via adversarial examples[C]//Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security. New York: ACM Press, 2019: 259-274.

[11] SHEJWALKAR V, HOUMANSADR A. Membership privacy for machine learning models through knowledge transfer[J]. Proceedings of the AAAI Conference on Artificial Intelligence, 2021, 35(11): 9549-9557.

[12] MAZZONE F, HEUVEL L D, HUBER M, et al. Repeated knowledge distillation with confidence masking to mitigate membership inference attacks[C]//Proceedings of the 15th ACM Workshop on Artificial Intelligence and Security. New York: ACM Press, 2022: 13-24.

[13] JARIN I, ESHETE B. MIAShield: defending membership inference attacks via preemptive exclusion of members[J]. Proceedings on Privacy Enhancing Technologies, 2023, 2023(1): 400-416.

[14] WANG X D, WU L F, GUAN Z T. GradDiff: gradient-based membership inference attacks against federated distillation with differential comparison[J]. Information Sciences, 2024, 658: 120068.

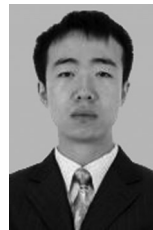
- [15] SHACHOR S, RAZINKOV N, GOLDSTEEN A, et al. Improving membership inference attacks against classification models[C]//International KES Conference on Intelligent Decision Technologies. Singapore: Springer Nature Singapore, 2025: 169-179.
- [16] HE X L, XU Y, ZHANG S C, et al. Enhance membership inference attacks in federated learning[J]. Computers & Security, 2024, 136: 103535.
- [17] CARLINI N, CHIEN S, NASR M, et al. Membership inference attacks from first principles[C]//Proceedings of the 2022 IEEE Symposium on Security and Privacy (SP). Piscataway: IEEE Press, 2022: 1897-1914.
- [18] SONG L W, MITTAL P. Systematic evaluation of privacy risks of machine learning models[C]//Proceedings of the USENIX Security Symposium. Berkeley: USENIX Association, 2025: 1-19.
- [19] WATSON L, GUO C, CORMODE G, et al. On the importance of difficulty calibration in membership inference attacks[J]. arXiv Preprint, arXiv: 2111.08440, 2021.
- [20] YE J Y, MADDI A, MURAKONDA S K, et al. Enhanced membership inference attacks against machine learning models[C]//Proceedings of the 2022 ACM SIGSAC Conference on Computer and Communications Security. New York: ACM Press, 2022: 3093-3106.
- [21] GRETTON A, BORGFARDT K M, RASCH M J, et al. A kernel two-sample test[J]. The Journal of Machine Learning Research, 2012, 13(1): 723-773.



王月琮 (1999-), 女, 甘肃兰州人, 西安电子科技大学博士生, 主要研究方向为通信及人工智能安全。



韩雪雪 (1994-), 女, 河北石家庄人, 西安电子科技大学博士生, 主要研究方向为网络安全。



张嘉伟 (1985-), 男, 山西太原人, 博士, 西安电子科技大学讲师、硕士生导师, 主要研究方向为数据安全、无线网络安全、AI安全、区块链和云安全。

#### [作者简介]



牛俊 (1992-), 女, 陕西西安人, 博士, 西安电子科技大学在站博士后, 主要研究方向为人工智能安全。



李兴华 (1978-), 男, 河南南阳人, 博士, 西安电子科技大学教授、博士生导师, 主要研究方向为无线网络安全与隐私保护、智能无人系统安全。



沈括 (1999-), 男, 宁夏中卫人, 西安电子科技大学硕士生, 主要研究方向为人工智能安全。



张玉清 (1966-), 男, 陕西宝鸡人, 博士, 中国科学院大学教授、博士生导师, 主要研究方向为网络与系统安全。