# Privacy Leaks by Adversaries: Adversarial Iterations for Membership Inference Attack

**Jing Xue***

**Zhishen Sun**    **Haishan Ye**    **Luo Luo**    **Xiangyu Chang**    **Ivor Tsang**    **Guang Dai**

## Abstract

Membership inference attack (MIA) has become one of the most widely used and effective methods for evaluating the privacy risks of machine learning models. These attacks aim to determine whether a specific sample is part of the model's training set by analyzing the model's output. While traditional membership inference attacks focus on leveraging the model's posterior output, such as confidence on the target sample, we propose `IMIA`, a novel attack strategy that utilizes the process of generating adversarial samples to infer membership. We propose to infer the member properties of the target sample using the number of iterations required to generate its adversarial sample. We conduct experiments across multiple models and datasets, and our results demonstrate that the number of iterations for generating an adversarial sample is a reliable feature for membership inference, achieving strong performance both in black-box and white-box attack scenarios. This work provides a new perspective for evaluating model privacy and highlights the potential of adversarial example-based features for privacy leakage assessment.

## 1   Introduction

Machine learning has widespread applications in many fields, such as autonomous driving [1, 33], medical [15] and financial systems [8, 23]. Training a model requires collecting a large amount of data and aims to help the model learn knowledge that generalizes well from the training data through multiple epochs. For example, a hospital may train a diagnostic model using thousands of patients' CT scans and treatment outcomes. While, this model is used to assist in diagnosis and treatment, several studies [7, 22, 24] have shown that neural network models tend to remember their training data and an adversary can exploit this weakness to launch membership inference attack (MIA) [6, 12, 25]. In such case, an adversary could infer whether a particular patient's record was used during the training process - potentially disclosing sensitive information like diagnosis results.

As a fundamental method to evaluate the privacy risk of machine learning models, membership inference attack (MIA) has received a lot of attention in recent years [2, 13, 24, 30]. Specifically, given a target model, an adversary aims to know if a target sample was part of the model's training set (being a member) or not (being a non-member). Successful membership inference attack can reveal individual's health, consumption habits, or even home location [5]. Therefore, studying attack methods such as MIA is important for understanding and evaluating the privacy risk of machine learning models.

---

*xuejing0203@stu.xjtu.edu.cn

Preprint. Under review.

Membership inference attacks (MIA) typically fall into two categories. Distribution-based MIA methods rely on differences between the training and test data distributions but require large datasets and shadow models, making them resource-intensive [6, 9, 29]. In contrast, metric-based MIA methods infer membership from the model's output, such as confidence scores, without access to the training data or shadow models [11, 26, 32]. However, these methods are limited to scenarios where the model exposes soft outputs (e.g., probabilities). For instance, the Softmax Response attack is effective only when the target model outputs confidence values and fails when only hard labels are available.

These limitations prompt us to ask: whether there exists a universal method that can solve these limitations, and remain effective in black-box as well as white-box scenarios without requiring extensive data or computation.

In this paper, we affirmatively answer this question by proposing a novel membership inference attack method, Iterations for Membership Inference Attack (IMIA), from the lens of adversarial samples' generation. Our key observation is that member samples, being closer to the decision boundary, generally require more iterations to generate adversarial examples than non-member samples. IMIA leverages this iteration gap across different settings—including white-box and black-box—by employing suitable adversarial attack strategies such as HopSkipJumpAttack, SimBA, and PGD [10, 17, 18]. Unlike prior work that relies on shadow models or access to similar data distributions [6, 29], IMIA operates without requiring the training set, making it lightweight and broadly applicable.

In general, our contributions are summarized as follows:

1. In this study, we propose a novel member inference attack method IMIA to infer whether the data belongs to the training set by analyzing the number of iterations required to generate adversarial samples. Different from traditional attack methods based on the posterior output of the target model, IMIA focuses on the generation process of adversarial samples, providing a tool to evaluate privacy leaks from the perspective of the internal operation of the model.

2. IMIA does *not* require any training data and training shadow models. The target sample is sufficient for IMIA to execute the attack. This strategy leverages the number of iterations required to generate adversarial samples from the target sample for MIA instead of posterior outputs of the target model.

3. IMIA is highly adaptable and universal. Our proposed method IMIA can be exploited in all settings compared with previous methods that can only be used in one specific situation. We have conducted experiments on multiple models and datasets, covering different network architectures and data distributions. Though our proposed method is simple, experimental results show that our method can effectively evaluate the privacy leakage risk of the model under both black-box and white-box settings.

## 2  Background and Related Work

In this section, we review research on membership inference attack, adversarial samples, and existing methods that we use as baselines.

### 2.1  Membership Inference Attack

Membership inference attack has achieved great attention because it revealed that machine learning models have serious risks of privacy leakage and remember its training data [2, 4, 20, 21, 27]. In membership inference attack, given the target sample $x$, the adversary aims to infer whether this sample is in the training set $D_{tr}$ of the target model $f_\theta$. As a result, membership inference attack can be seen as a binary classification privacy game. The participants in this game are the challenger $\mathcal{C}$ and the adversary $\mathcal{A}$. The game process can be broken down into several steps:

1. The challenger samples a training set $D_{tr} \sim \mathbb{D}$, and trains the target model $f_\theta$.

2. The challenger randomly chooses a bit $b \in \{0, 1\}$. If $b = 0$, he will choose the target sample $(x, y) \in \mathbb{D}$ where $y$ is the ground-truth label of the target sample $x$, but $(x, y)$ is not in $D_{tr}$. If $b = 1$, he will choose the target sample in $D_{tr}$ directly.

3. The challenger sends the target sample to the adversary and allows the adversary to query the target model.

4. The adversary gets the target sample and returns a bit $\hat{b}$ by querying the target model. If $\hat{b} = b$, the adversary will win this game.

Based on the adversary's ability to access the target model, MIA can be divided into two categories:

**Black-box Membership Inference.** In the black-box setting, the adversary can only access the posterior output of the model (like confidence or hard labels) [19, 31]. This is a true scenario in the real world. In the black-box case, the attack methods can also be divided into two categories: the first is a distribution-based MIA [6, 9, 13, 29], in which the adversary needs to train a lot of additional shadow models as a proxy to mimic the target model. These shadow models require a large amount of data in the same distribution as the target sample and use the same model framework. The adversary trains shadow models with and without the target sample to learn the distribution difference between members and non-member samples, and then judges the member attributes of the target sample.

The second is the metric-based MIA, in which the adversary does not need to train additional shadow models. In Step 4, the adversary only uses the posterior output of the target model and designs a metric $\mathrm{Me}(\cdot)$, such as Softmax Response [19, 26], Prediction Entropy [22, 32] and Modified Entropy [25]. Specifically, Softmax Response [26] computes the output probabilities of the target model $f_\theta(x)$ for the input sample $x$, obtaining the predicted probability $f_\theta(x)_i$ for each class $i$, and compares it with a preset threshold $\tau$. If the maximum predicted probability exceeds the threshold, the adversary infers that the sample belongs to the training set. Formally,

$$I_{soft}(f_\theta, (x, y)) = \mathbb{1}\{\max f_\theta(x)_i \geq \tau\}$$

Salem et al. [22] proposed to use the Prediction Entropy to conduct MIA. Prediction Entropy measures the uncertainty of the model's prediction and they thought that the prediction entropy of the member data would be close to 0 and a larger entropy value to the non-member data on the target model. This is formally expressed as:

$$I_{ent}(f_\theta, (x, y)) = \mathbb{1}\{-\sum_i f_\theta(x)_i \log(f_\theta(x)_i) \leq \tau\}$$

Besides, Song and Mittal [25] proposed Modified Entropy that considered the ground-truth label of the target sample by decreasing the uncertainty on the true label of the target sample and increasing the uncertainty on the wrong label. It can be computed through:

$$Mentr(f_\theta(x), y) = -(1 - f_\theta(x)_y)\log(f_\theta(x)_y)$$
$$-\sum_{i \neq y} f_\theta(x)_i \log(1 - f_\theta(x)_i)$$

They inferred a target sample in $\mathrm{D}_{tr}$ according to:

$$I_{Mentr}(f_\theta, (x, y)) = \mathbb{1}\{Mentr f_\theta(x, y) \leq \tau\}$$

That is if the modified entropy value of the target sample is lower than the preset threshold, it will be recognized as a member.

**White-box Membership Inference.** In the white-box setting, the adversary can not only get the posterior output of the target model but also the internal parameters of the target model, like the loss and gradient information during the progress of the target model training [3, 16]. Yeom et al. [32] used the internal weights of the model and the loss of the target sample on the model to conduct MIA. If an adversary can get access to the logits output of the target model, he can conduct MIA through:

$$I_{loss}(f_\theta, (x, y)) = \mathbb{1}\{-l(f_\theta(x), y) \leq \tau\}$$

where the loss function will be cross-entropy loss. The sample $x$ which has lower loss value in the training set $\mathrm{D}_{tr}$. This capability allows an adversary to obtain more sensitive information of the victim model, but this capability is rarely available in the real world.

## 2.2 Adversarial sample

In a deep neural network, adversarial samples are inputs intentionally perturbed with small but deliberate perturbations that can cause deep neural networks to make incorrect predictions. Goodfellow et al. [16] first revealed the vulnerability of neural networks to such inputs. Then, numerous attack methods have been proposed, including white-box approaches like PGD [18] and CW [3], as well as black-box approaches like SimBA [17], which provide effective strategies for generating adversarial examples.

In the background of MIA, the characteristics of adversarial samples are used to infer the privacy of the model's training data. Song et al. [26] used the confidence output on adversarial samples to judge the member attributes of the target samples. They believed that due to the robustness of the model, the adversarial samples generated from the member data show relatively stable prediction results on the model. Afterwards, Choquette-Choo et al. [12] found that the distance from the adversarial sample to the model's decision boundary can be directly used to carry out MIA and had a great performance under the black-box setting. They judged a sample as a member if the distance from the adversarial sample to its decision boundary is larger than a preset threshold:

$$I_{dis}(f_\theta, (x, y)) = \mathbb{1}\{d(x, \hat{x}) \geq \tau\}$$

Del Grosso et al. [14] also analyzed MIA from this perspective.

In our paper, we carry out the metric-based black-box membership inference from the lens of the generation of adversarial samples. Different from prior works[12], we do not measure the boundary distance between the adversarial sample and its decision boundary but record the number of iterations during the process of generating adversarial samples. Softmax Response, Prediction Entropy, and Modified Entropy as typical representatives in the score-based metric attacks, we choose these as our baselines. In the case where the target model only outputs hard labels, we choose the boundary distance attack as our baseline because boundary distance has a great performance in the metric-based condition.
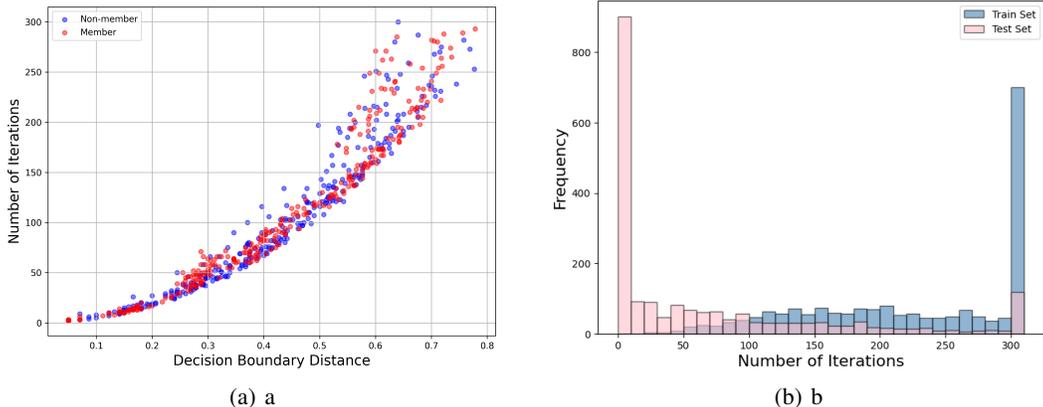


(a) a          (b) b

Figure 1: (a) Scatter diagram showing the relationship between the distance from samples to the decision boundary and the number of iterations required to generate adversarial examples using SimBA for ResNet trained on CIFAR10. (b) Histogram about the number of iterations per-sample over 2k samples from the training set (blue) and the same number from the testing set. The adversarial samples are generated using SimBA for ResNet trained on CIFAR100.

## 3 Membership Inference Attack during Adversarial Examples Iteration

In this section, we discuss the methodology of this paper, which aims to reveal privacy risks of the target model from the lens of adversarial samples' generation process.

4

## 3.1 Motivation

While prior MIA methods often rely on confidence scores or boundary distance [12, 26], we find that boundary distance may not reliably distinguish members from non-members because different samples can have similar distance regardless of membership, as shown in Figure 1(a). We plotted a scatter diagram showing the relationship between the distance from samples to the decision boundary and the number of iterations required to generate adversarial samples.

Instead, we observe that the member samples generally require more iterations to generate adversarial samples than non-members. We use "SimBA" [17] to generate adversarial samples and record the number of iterations to generate its adversarial sample for each sample. The histogram is shown in Figure 1(b): the samples from the training set (blue, member) need more iterations than those from the test set (pink, non-member).

This insight reveals a new, consistent signal for membership inference and motivates our method, IMIA, which leverages the number of iterations required to generate an adversarial sample for a target sample. If this number exceeds a preset threshold, this target sample will be inferred as a member; otherwise, as a non-member. The overall framework of `IMIA` is shown in Figure 2.
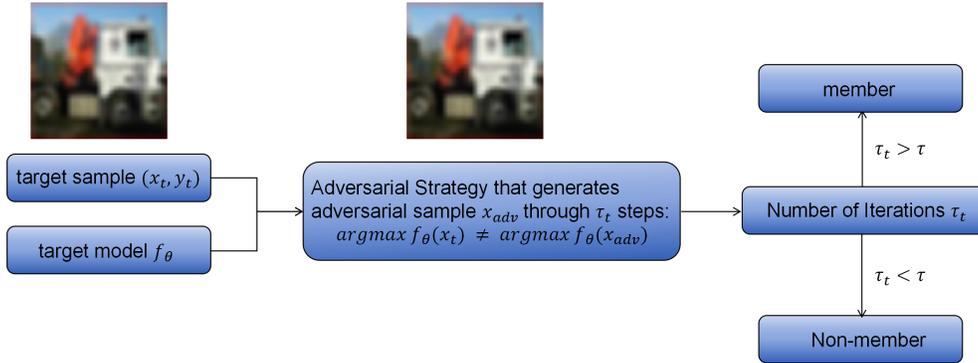


Figure 2: Diagrammatic sketch for `IMIA` to conduct MIA.

## 3.2 Methodology

Given the target model and the target images, the adversary can choose an adversarial strategy $\mathcal{S}$ in SimBA [17], HopSkipJumpAttack [10] and PGD [18] based on different MIA settings to generate adversarial samples and measure the number of iterations during this process.

**Score-based black-box attacks.** Adversary can obtain the full probability output of the target sample. As a result, we choose SimBA [17] to conduct adversarial attack which provides a simple but efficient strategy to change the target sample's output. The optimization goal in SimBA is to minimize the probability on the true label $y$ of the target sample $x$ :

$$\min_{\delta} \quad p_{f_\theta}(y|x + \delta)$$

$$\text{subject to:} \quad \|\delta\|_2 < d, \text{ queries} \leq M,$$

where $\delta$ represents the perturbations and $M$ is the budget for the number of queries to the target model. SimBA solves this problem through randomly selecting a predefined orthonormal basis and either adding or subtracting it from the target sample according to the confidence scores which are checked if the target sample moves toward the decision boundary. During the process of generating adversarial samples for the target sample, we are concerned about the number of iterations after getting adversarial samples successfully.

**Decision-based black-box attacks.** Different from score-based black-box attacks, the adversary can only obtain the label of the target input without any other information. In this case, we choose

---
**Algorithm 1** IMIA
---
**Require:** Target model $f_\theta$, target sample $(x_t, y_t)$, adversarial strategy $\mathcal{S}$, threshold $\tau$
  1: $\tau_t \leftarrow \mathcal{S}(f_\theta, (x_t, y_t))$
**Ensure:** Number of iterations $\tau_t$ for generating adversarial sample
  2: **if** $\tau_t \geq \tau$ **then**
  3:     **return** $\mathbb{1}(\tau_t \geq \tau)$                                 ▷ Classify as member
  4: **else**
  5:     **return** $0$                                        ▷ Classify as non-member
  6: **end if**
---

"HopSkipJumpAttack" [10] whose performance is close to the white-box attack. Given the target image $(x, y)$, adversary starts from randomly choosing a point $x'$ that is not classified to label $y$ by the target model and walks along the decision boundary to minimize the distance between the original image $x$ to its adversarial image $x'$. In our methodology, we measure the number of iterations that adversarial samples can satisfy our request.

**White-box attacks.** In the white-box setting, We choose Projected Gradient Descent (PGD) [18] to generate adversarial samples. Formally, adversarial sample $x_{adv}$ is generated by:

$$x_{adv} = \text{Clip}_{x,\epsilon}(x_N^{adv} + \alpha \, \text{sign}\left(\nabla_x J(x_{N+1}^{adv}, y)\right)$$

In `IMIA`, we will choose one of them to generate adversarial samples according to different MIA settings. The pseudocode for `IMIA` is listed in Algorithm 1. As shown in Algorithm 1, given a target model $f_\theta$, a target sample $(x_t, y_t)$, an adversarial attack strategy $S$, and a threshold $\tau$, IMIA first computes the number of iterations $\tau_t$ required by $S$ to generate a successful adversarial example. The adversarial strategy $S$ is chosen based on the attack setting including PGD, SimBA and HopSkipJumpAttack as we described before. If $\tau_t \geq \tau$, the sample is classified as a *member*; otherwise, it is classified as a *non-member*.

In this paper, we consider three different attack strategies for different MIA settings, and we demonstrate the implementation details for each adversarial strategy in Appendix A. Note that we do not propose a novel method for adversarial attacks, instead, we only care about how the adversarial sample is generated from the original target sample.

# 4 Evaluation

In this section, we will evaluate the effectiveness and universality of our algorithm.

## 4.1 Experiment Setup

**Datasets.** In our experiment, we consider three different datasets which are all common in image recognition tasks. **CIFAR10** and **CIFAR100** all include 60k images and can be split into 10 and 100 classes respectively. In PyTorch, CIFAR10 and CIFAR100 are divided into 50k images in the training dataset and others are in the test dataset. **STL-10** includes 5k images in the training set and 8k images in its test set. In our experiment, we use all samples from the training dataset to train the target model for each model architecture, and samples in the test set which do not participate in the training process are used to validate the model's accuracy.

**Target Model.** In our experiments, we use four different model architectures: ResNet50, VGG19, ResNeXt29_2x64d, and DenseNet121. These models represent a diverse set of machine learning architectures. We train each model for 100 epochs in every dataset during the training process to ensure that the models are sufficiently trained and can produce meaningful outputs for membership inference attack. And we use the test set of the target model to verify the accuracy of the target model avoiding the low accuracy of the target model on the test set.

**Metrics.** We use a balanced evaluation set to evaluate our method and report the inference accuracy, AUROC scores, and the false positive rate (fpr) under different true positive rate (tpr). The inference accuracy considers both the true positive rate and the false positive rate and gives 50% if the adversary guesses randomly. Area Under the Receiver Operating Characteristic Curve (AUROC) is the area

under ROC curve, which is obtained by plotting the ratio of TPR to FPR at different thresholds. The closer the AUROC value is to 1, the better the attack performance.

In a balanced evaluation set, there is an equal number of member samples and non-member samples of the target model. As a result, our evaluation set consists of 3k samples from the target model's training set as member samples and uniformly selects the same number of samples from the test set as non-member samples. We repeat this many times to get different combinations of evaluation sets and finally report results on average.

In our paper, IMIA is a metric-based universal attack method which does not need any training data or shadow models, so to ensure a fair comparison, we also choose the metric-based method as our baselines instead of these methods using shadow models. The specific reasons can be found in the Appendix B. Moreover, we display the specific time cost of computation to run IMIA and compare this to the cost of these methods based on shadow models in Appendix C.

Table 1: Membership inference results for different score-based attacks on CIFAR10, CIFAR100, and STL10 datasets in the black-box settings. We choose "SimBA" to generate adversarial samples for measuring the number of iterations for target samples. The evaluation set is composed of 3k samples from the training set and 3k samples from the test set of the target model, and we repeat this procedure 20 times. The inference accuracy(%) and AUROC(%) are reported in the average value with standard deviation. The best results in each case are in bold.

| Strategy | ResNet | | ResNeXt | | VGG | | DenseNet | |
|---|---|---|---|---|---|---|---|---|
| | AUROC ↑ | Accuracy ↑ | AUROC | Accuracy | AUROC | Accuracy | AUROC | Accuracy |
| **CIFAR100** | | | | | | | | |
| Softmax | 88.91± 0.13 | 84.57± 0.15 | 63.66± 0.09 | 60.42± 0.10 | 63.73± 0.12 | 59.94± 0.18 | 67.70± 0.11 | 63.26± 0.14 |
| Entropy | 88.97± 0.08 | 84.59± 0.11 | 64.22± 0.13 | 54.42± 0.07 | 64.95± 0.14 | 61.26± 0.12 | 68.78± 0.16 | 64.00± 0.11 |
| Mentr. | **89.24± 0.15** | **85.43± 0.16** | 68.42± 0.10 | 64.48± 0.12 | **69.61± 0.14** | **65.41± 0.13** | 73.19± 0.12 | 67.90± 0.18 |
| IMIA(ours) | 89.11± 0.29 | 83.62± 0.29 | **68.93± 0.58** | **64.82± 0.55** | 67.36± 0.64 | 64.68± 0.55 | **73.80± 0.48** | **68.91± 0.40** |
| **CIFAR10** | | | | | | | | |
| Softmax | 70.97± 0.05 | 65.84± 0.15 | 74.02± 0.13 | 68.05± 0.19 | 64.15± 0.16 | 60.60± 0.15 | 69.21± 0.15 | 64.17± 0.14 |
| Entropy | 71.05± 0.13 | 65.87± 0.29 | **74.35± 0.14** | **68.24± 0.17** | 64.39± 0.15 | 60.70± 0.12 | 69.51± 0.16 | 64.27± 0.14 |
| Mentr. | 71.95± 0.12 | 66.17± 0.16 | 73.54± 0.15 | 67.95± 0.11 | 64.56± 0.12 | 60.58± 0.14 | 69.44± 0.15 | 64.23± 0.12 |
| IMIA(ours) | **74.43± 0.15** | **68.35± 0.13** | 73.34± 0.09 | 67.77± 0.15 | **66.53± 0.13** | **62.87± 0.15** | **75.20± 0.09** | **68.97± 0.15** |
| **STL10** | | | | | | | | |
| Softmax | 56.49± 0.13 | 55.53± 0.15 | 61.68± 0.17 | 59.43± 0.11 | 57.06± 0.16 | 56.63± 0.12 | 72.97± 0.12 | 68.52± 0.13 |
| Entropy | 55.91± 0.18 | 55.28± 0.19 | 62.08± 0.12 | 59.72± 0.11 | 57.75± 0.14 | 57.61± 0.13 | 73.22± 0.13 | 68.59± 0.15 |
| Mentr. | 61.05± 0.11 | 59.02± 0.13 | **70.84± 0.18** | **66.37± 0.10** | **62.68± 0.14** | **59.92± 0.12** | 78.41± 0.11 | 73.97± 0.09 |
| IMIA(ours) | **61.26± 0.15** | **59.70± 0.12** | 69.25± 0.15 | 66.35± 0.14 | 61.32± 0.08 | 59.87± 0.14 | **81.18± 0.12** | **75.41± 0.11** |

## 4.2 Results

We analyze our results in the black-box and the white-box settings separately. In the black-box setting, there are two types: one is score-based attack where the target model outputs the confidence and labels at the same time; the other is the target model that only outputs hard labels. In the white-box setting, adversaries have access to the loss during the process of generating adversarial samples. Specifically, IMIA resorts to generate adversarial samples like "PGD" [18, 28] for white-box MIA attack, "SimBA" [17] for black-box MIA attack and "HopSkipJumpAttack" for hard labels MIA attack.

### 4.2.1 score-based setting

In the score-based setting, the target model outputs the confidence and labels at the same time. In this case, we choose "SimBA" [17] as our strategy to generate corresponding adversarial samples of the original samples. Softmax Response [26], Prediction Entropy [22] and Modified Entropy [25] as typical representatives in the score-based metric membership inference attack, we choose these as our baselines.

Table 1 shows the results of our attacks and comparisons between IMIA and other baselines in the score-based setting. The results show our strategy that depending on the number of iterations performs well in distinguishing member samples and non-member samples. For example, for our proposed attack against CIFAR10 DenseNet classifier, the membership inference AUROC is increased from 69.21 to 75.20 on average and the accuracy is increased from 64.17 to 68.97. In other

Table 2: Membership inference results for different decision-based attacks that only output hard labels on CIFAR10, CIFAR100 and STL10 datasets. We choose "HopSkipJumpAttack" strategy, and measure the distance and the number of iterations between original samples and adversarial samples. The inference accuracy(%) and AUROC(%) are reported in the average value with standard deviation. The best results in each case are in bold.

| Strategy | ResNet | | ResNeXt | | VGG | | DenseNet | |
|---|---|---|---|---|---|---|---|---|
| | AUROC ↑ | Accuracy ↑ | AUROC | Accuracy | AUROC | Accuracy | AUROC | Accuracy |
| **CIFAR100** | | | | | | | | |
| Boundary | **86.12±0.35** | **81.59±0.43** | **69.49±0.29** | **66.90±0.34** | 63.29±0.33 | 61.90±0.26 | 69.72±0.32 | 66.90±0.33 |
| IMIA(ours) | 85.95±0.27 | 81.39±0.25 | 69.20±0.21 | 66.70±0.25 | **63.39±0.31** | **62.39±0.30** | **70.03±0.22** | **67.20±0.22** |
| **CIFAR10** | | | | | | | | |
| Boundary | 72.84±0.30 | 66.70±0.25 | **70.43±0.26** | **65.60±0.31** | 66.50±0.23 | 63.0±0.21 | 69.72±0.15 | 66.00±0.22 |
| IMIA(ours) | **73.12±0.25** | **67.39±0.27** | 69.77±0.25 | 65.10±0.23 | **69.31±0.21** | **65.10±0.21** | **71.81±0.19** | **66.90±0.17** |
| **STL10** | | | | | | | | |
| Boundary | 61.45±0.27 | 60.24±0.31 | 70.82±0.35 | 67.34±0.31 | 63.78±0.24 | 61.60±0.20 | 82.42±0.13 | 75.29±0.17 |
| IMIA(ours) | **62.00±0.19** | **60.41±0.22** | **72.21±0.25** | **69.38±0.20** | **63.93±0.18** | **61.63±0.24** | **83.36±0.19** | **76.57±0.21** |

Table 3: Membership inference results for Loss and our method on CIFAR10, CIFAR100, and STL10 in the white-box setting. We choose "PGD" to generate adversarial samples and for the loss, we choose the cross entropy loss for each target sample. The inference accuracy(%) and AUROC(%) are reported in the average value with standard deviation. The best results in each case are in bold.

| Strategy | ResNet | | ResNeXt | | VGG | | DenseNet | |
|---|---|---|---|---|---|---|---|---|
| | AUROC ↑ | Accuracy ↑ | AUROC | Accuracy | AUROC | Accuracy | AUROC | Accuracy |
| **CIFAR100** | | | | | | | | |
| Loss | 89.24± 0.24 | 85.68± 0.26 | 68.47± 0.46 | 65.14± 0.51 | 69.95± 0.59 | 65.87± 0.60 | 73.31± 0.43 | 68.42± 0.46 |
| Boundary | 88.37± 0.31 | 82.82± 0.26 | 66.16±0.45 | 63.97±0.49 | 64.69±0.50 | 62.62±0.51 | 66.43 ±0.40 | 64.46±0.43 |
| IMIA(ours) | **96.12± 0.21** | **90.54± 0.28** | **69.31± 0.54** | **65.82± 0.51** | **71.47± 0.45** | **68.55± 0.42** | **73.41± 0.40** | **68.91± 0.39** |
| **CIFAR10** | | | | | | | | |
| Loss | 72.26± 0.49 | 66.61± 0.45 | 74.33± 0.35 | 68.96± 0.39 | 65.09± 0.36 | 61.72± 0.44 | 69.94± 0.35 | 64.79± 0.32 |
| Boundary | **75.45±0.47** | **69.03±0.45** | 74.19±0.37 | 68.62±0.37 | 65.94±0.40 | 61.67±0.42 | 75.69±0.33 | 69.82±0.29 |
| IMIA(ours) | 74.49± 0.46 | 68.97± 0.39 | **74.45± 0.34** | **69.83± 0.31** | **66.96± 0.35** | **62.67± 0.32** | **76.29± 0.33** | **70.15± 0.31** |
| **STL10** | | | | | | | | |
| Loss | 61.52± 0.45 | 59.20± 0.33 | **71.01± 0.45** | **66.39± 0.42** | **62.94± 0.49** | **59.68± 0.36** | 78.57± 0.46 | 74.06± 0.45 |
| Boundary | 60.29±0.43 | 58.75±0.32 | 68.38±0.50 | 65.25±0.47 | 59.66±0.45 | 59.05±0.35 | 81.78±0.46 | 75.63±0.47 |
| IMIA(ours) | **61.94± 0.47** | **59.64± 0.36** | 70.07± 0.53 | 66.19± 0.51 | 61.46± 0.52 | 59.61± 0.49 | **82.19± 0.61** | **75.63± 0.50** |

words, our strategy can effectively reveal the privacy risk of the target model during the process of generating adversarial samples. Figure 3 show the false positive rate under different true positive rate corresponding to Table 1. The figures illustrate how the false positive rate varies as the true positive rate changes. Our proposed attack has a relatively lower fpr than others.

### 4.2.2 decision-based setting

In another black-box setting called decision-based attack, the target model only outputs hard labels. In this case, we use "HopSkipJumpAttack" [10] strategy to generate adversarial samples. We then compare our attack with the " Decision boundary" method [12] which has strong performance when the target model only outputs hard labels. The "Boundary" measures the distance from the adversarial samples to their decision boundary. Our `IMIA` measures the number of iterations during the generation of adversarial samples. The comparison results are shown in Table 2. We can observe that `IMIA` has advantages over "Boundary" method. For example, the membership inference AUROC and accuracy in STL10 are all higher than "Boundary" method for all classifiers. Furthermore, for the VGG model on the CIFAR10, our `IMIA` achieves an accuracy about $2.8\%$ higher than "Boundary". Though `IMIA` is simple, in this most difficult situation, it can still work efficiently.

### 4.2.3 white-box setting

We also evaluate `IMIA` in the white-box setting and show the comparison in Table 3 where the best results in each case are in bold. In the white-box setting, "Loss" method [32] is one of the most common methods to measure the degree of privacy leaks, so we also use "Loss" method as baseline

and compute the cross entropy loss to quantify the privacy risks associated with the target model. At the same time, we also compare the results of "Decision Boundary" method. In this setting, we use "PGD" [18] methodology to generate adversarial samples. Table 3 shows that `IMIA` surpasses the "Loss" method both in the inference accuracy and AUROC in most classifiers and datasets. Especially, when applied to the DenseNet architecture on CIFAR10 dataset, `IMIA` can achieve an accuracy over 5% higher than the "Loss" method.



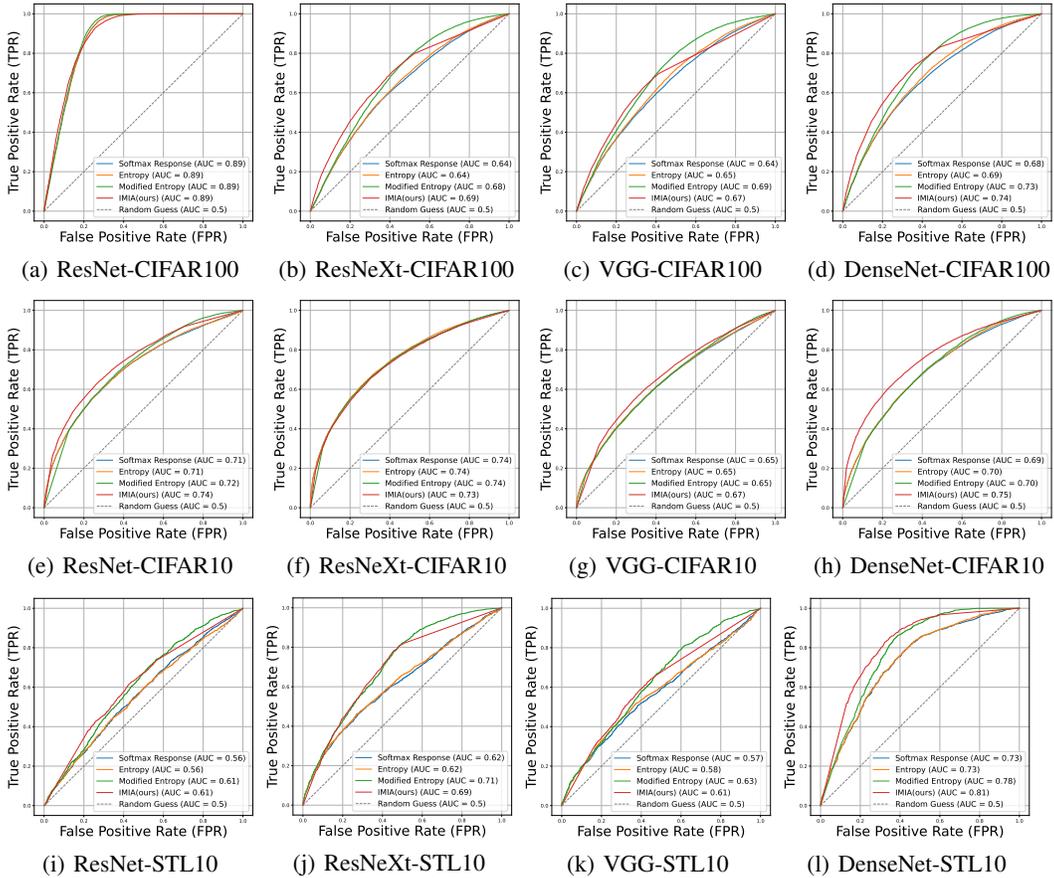|  |  |  |  |
|---|---|---|---|
| (a) ResNet-CIFAR100 | (b) ResNeXt-CIFAR100 | (c) VGG-CIFAR100 | (d) DenseNet-CIFAR100 |
| (e) ResNet-CIFAR10 | (f) ResNeXt-CIFAR10 | (g) VGG-CIFAR10 | (h) DenseNet-CIFAR10 |
| (i) ResNet-STL10 | (j) ResNeXt-STL10 | (k) VGG-STL10 | (l) DenseNet-STL10 |

Figure 3: ROC curve on MIA for the combination of different models on CIFAR100, CIFAR10 and STL-10. They are drawn on the balanced evaluation set and correspond to Table 1.

### 4.3 Summary

In three different application conditions, the increasing inference accuracy and AUROC value prove that our methodology `IMIA` has great performance in distinguishing the member data and the non-member data. Even in the most difficult situation, our methodology can still work well. All results show that `IMIA` has great adaptive ability and can be applied in both white-box and black-box settings without knowing data from the training set.

## 5 Conclusion

In this paper, We propose a universal membership inference attack method, called `IMIA`, which performs the membership inference attack from the perspective of adversarial samples' generation process. The key idea of `IMIA` is to measure the number of iterations by the process of generating the adversarial samples, and use this metric to infer whether the target samples belong to the model's training set or not. `IMIA` with different adversarial strategies can be applied in different settings. Accordingly, we conduct experiments in different MIA settings and on different datasets such as CIFAR10, CIFAR100 and STL10 datasets under different model architectures. Our experiment

results show that our proposed methodology has great performance in different situations with higher AUROC values and inference accuracy compared to the other metric-based MIA algorithms. All experiments highlight the superior performance of `IMIA` and prove that `IMIA` is universal and adaptable in most settings to detect the privacy risk of the target model while requiring fewer computational resources, making it a more efficient choice for MIA.

# References

[1] Shunsuke Aoki, Issei Yamamoto, Daiki Shiotsuka, Yuichi Inoue, Kento Tokuhiro, and Keita Miwa. Superdriverai: Towards design and implementation for end-to-end learning-based autonomous driving. In *2023 IEEE Vehicular Networking Conference (VNC)*, pages 195–198. IEEE, 2023.

[2] Martin Bertran, Shuai Tang, Aaron Roth, Michael Kearns, Jamie H Morgenstern, and Steven Z Wu. Scalable membership inference attacks via quantile regression. *Advances in Neural Information Processing Systems*, 36, 2024.

[3] Nicholas Carlini and David Wagner. Towards evaluating the robustness of neural networks. In *2017 ieee symposium on security and privacy (sp)*, pages 39–57. Ieee, 2017.

[4] Nicholas Carlini, Chang Liu, Úlfar Erlingsson, Jernej Kos, and Dawn Song. The secret sharer: Evaluating and testing unintended memorization in neural networks. In *28th USENIX security symposium (USENIX security 19)*, pages 267–284, 2019.

[5] Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar Erlingsson, et al. Extracting training data from large language models. In *30th USENIX Security Symposium (USENIX Security 21)*, pages 2633–2650, 2021.

[6] Nicholas Carlini, Steve Chien, Milad Nasr, Shuang Song, Andreas Terzis, and Florian Tramer. Membership inference attacks from first principles. In *2022 IEEE Symposium on Security and Privacy (SP)*, pages 1897–1914. IEEE, 2022.

[7] Nicholas Carlini, Daphne Ippolito, Matthew Jagielski, Katherine Lee, Florian Tramer, and Chiyuan Zhang. Quantifying memorization across neural language models. *arXiv preprint arXiv:2202.07646*, 2022.

[8] Sotirios P Chatzis, Vassilis Siakoulis, Anastasios Petropoulos, Evangelos Stavroulakis, and Nikos Vlachogiannakis. Forecasting stock market crisis events using deep and statistical machine learning techniques. *Expert systems with applications*, 112:353–371, 2018.

[9] Harsh Chaudhari, Giorgio Severi, Alina Oprea, and Jonathan Ullman. Chameleon: Increasing label-only membership leakage with adaptive poisoning. *arXiv preprint arXiv:2310.03838*, 2023.

[10] Jianbo Chen, Michael I Jordan, and Martin J Wainwright. Hopskipjumpattack: A query-efficient decision-based attack. In *2020 ieee symposium on security and privacy (sp)*, pages 1277–1294. IEEE, 2020.

[11] Yufei Chen, Chao Shen, Yun Shen, Cong Wang, and Yang Zhang. Amplifying membership exposure via data poisoning. *Advances in Neural Information Processing Systems*, 35:29830–29844, 2022.

[12] Christopher A Choquette-Choo, Florian Tramer, Nicholas Carlini, and Nicolas Papernot. Label-only membership inference attacks. In *International conference on machine learning*, pages 1964–1974. PMLR, 2021.

[13] Edoardo Debenedetti, Giorgio Severi, Nicholas Carlini, Christopher A Choquette-Choo, Matthew Jagielski, Milad Nasr, Eric Wallace, and Florian Tramèr. Privacy side channels in machine learning systems. In *33rd USENIX Security Symposium (USENIX Security 24)*, pages 6861–6848, 2024.

[14] Ganesh Del Grosso, Hamid Jalalzai, Georg Pichler, Catuscia Palamidessi, and Pablo Piantanida. Leveraging adversarial examples to quantify membership information leakage. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10399–10409, 2022.

[15] Rohit R Dixit. Risk assessment for hospital readmissions: Insights from machine learning algorithms. *Sage Science Review of Applied Machine Learning*, 4(2):1–15, 2021.

[16] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014.

[17] Chuan Guo, Jacob Gardner, Yurong You, Andrew Gordon Wilson, and Kilian Weinberger. Simple black-box adversarial attacks. In *International conference on machine learning*, pages 2484–2493. PMLR, 2019.

[18] Aleksander Madry. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017.

[19] Milad Nasr, Reza Shokri, and Amir Houmansadr. Comprehensive privacy analysis of deep learning: Passive and active white-box inference attacks against centralized and federated learning. In *2019 IEEE symposium on security and privacy (SP)*, pages 739–753. IEEE, 2019.

[20] Milad Nasr, Nicholas Carlini, Jonathan Hayase, Matthew Jagielski, A Feder Cooper, Daphne Ippolito, Christopher A Choquette-Choo, Eric Wallace, Florian Tramèr, and Katherine Lee. Scalable extraction of training data from (production) language models. *arXiv preprint arXiv:2311.17035*, 2023.

[21] USVSN Sai Prashanth, Alvin Deng, Kyle O'Brien, Jyothir SV, Mohammad Aflah Khan, Jaydeep Borkar, Christopher A Choquette-Choo, Jacob Ray Fuehne, Stella Biderman, Tracy Ke, et al. Recite, reconstruct, recollect: Memorization in lms as a multifaceted phenomenon. *arXiv preprint arXiv:2406.17746*, 2024.

[22] Ahmed Salem, Yang Zhang, Mathias Humbert, Pascal Berrang, Mario Fritz, and Michael Backes. Ml-leaks: Model and data independent membership inference attacks and defenses on machine learning models. *arXiv preprint arXiv:1806.01246*, 2018.

[23] Aristeidis Samitas, Elias Kampouris, and Dimitris Kenourgios. Machine learning as an early warning system to predict financial crisis. *International Review of Financial Analysis*, 71: 101507, 2020.

[24] Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership inference attacks against machine learning models. In *2017 IEEE symposium on security and privacy (SP)*, pages 3–18. IEEE, 2017.

[25] Liwei Song and Prateek Mittal. Systematic evaluation of privacy risks of machine learning models. In *30th USENIX Security Symposium (USENIX Security 21)*, pages 2615–2632, 2021.

[26] Liwei Song, Reza Shokri, and Prateek Mittal. Privacy risks of securing machine learning models against adversarial examples. In *Proceedings of the 2019 ACM SIGSAC conference on computer and communications security*, pages 241–257, 2019.

[27] Marlon Tobaben, Gauri Pradhan, Yuan He, Joonas Jälkö, and Antti Honkela. Understanding practical membership privacy of deep learning. *arXiv preprint arXiv:2402.06674*, 2024.

[28] Florian Tramèr, Nicolas Papernot, Ian Goodfellow, Dan Boneh, and Patrick McDaniel. The space of transferable adversarial examples. *arXiv preprint arXiv:1704.03453*, 2017.

[29] Florian Tramèr, Reza Shokri, Ayrton San Joaquin, Hoang Le, Matthew Jagielski, Sanghyun Hong, and Nicholas Carlini. Truth serum: Poisoning machine learning models to reveal their secrets. In *Proceedings of the 2022 ACM SIGSAC Conference on Computer and Communications Security*, pages 2779–2792, 2022.

[30] Lauren Watson, Chuan Guo, Graham Cormode, and Alex Sablayrolles. On the importance of difficulty calibration in membership inference attacks. *arXiv preprint arXiv:2111.08440*, 2021.

[31] Yutong Wu, Han Qiu, Shangwei Guo, Jiwei Li, and Tianwei Zhang. You only query once: An efficient label-only membership inference attack. In *The Twelfth International Conference on Learning Representations*.

[32] Samuel Yeom, Irene Giacomelli, Matt Fredrikson, and Somesh Jha. Privacy risk in machine learning: Analyzing the connection to overfitting. In *2018 IEEE 31st computer security foundations symposium (CSF)*, pages 268–282. IEEE, 2018.

[33] Hanyi Yu, Shuning Huo, Mengran Zhu, Yulu Gong, and Yafei Xiang. Machine learning-based vehicle intention trajectory recognition and prediction for autonomous driving. In *2024 7th International Conference on Advanced Algorithms and Control Engineering (ICAACE)*, pages 771–775. IEEE, 2024.

[34] Sajjad Zarifzadeh, Philippe Liu, and Reza Shokri. Low-cost high-power membership inference attacks. *arXiv preprint arXiv:2312.03262*, 2023.

# A Implementation details for each adversarial strategy

In Table 4, we show the implementation details for the Simple Black-box Attack (SimBA) under the score-based setting. The parameters in our experiments are similar with Guo et al. [17]. The most difference is the perturbation size ($\epsilon$) which we set a small value for observing the number of iterations during the adversarial samples' generation process.

We implement the HopSkipJump Attack (HSJA) [10] under a decision-based black-box setting. We follow the standard setup of the HopSkipJump Attack (HSJA) [10], using $L_2$ constraint with input clipping to $[0, 1]$, and set the number of iterations to 100. Other hyperparameters, such as step size search strategy and evaluation limits, are set according to default values or tuned empirically for stability in Table 5.

We employ the Projected Gradient Descent (PGD) attack [18] under the white-box setting with an $L_\infty$ perturbation budget of $\epsilon = \frac{3}{255}$, step size $\alpha = 0.001$, and 50 attack steps. A random start is enabled to initialize the attack from a perturbed point within the allowed norm ball. The attack is untargeted.

Table 4: SimBA Attack Hyperparameters

| Parameter | Value |
|---|---|
| Maximum iterations (`max_iters`) | 300 |
| Frequency dimensions (`freq_dims`) | 32 |
| Stride (`stride`) | 7 |
| Perturbation size ($\epsilon$) | 0.05 |
| $L_\infty$ bound (`linf_bound`) | 0.0 (unbounded) |
| Perturbation order (`order`) | `rand` |
| Targeted attack (`targeted`) | False |
| Pixel attack (`pixel_attack`) | False (frequency domain) |
| Logging interval (`log_every`) | 40 |

Table 5: HSJA Attack Hyperparameters

| Parameter | Value |
|---|---|
| Clipping Range (`clip_min`, `clip_max`) | [0, 1] |
| Norm Constraint (`constraint`) | $L_2$ |
| Number of Iterations (`num_iterations`) | 100 |
| Binary Search Confidence (`gamma`) | 1.0 |
| Step Size Search Strategy (`stepsize_search`) | Geometric progression |
| Max Gradient Evaluations (`max_num_evals`) | 1e4 |
| Initial Gradient Evaluations (`init_num_evals`) | 100 |

# B Justification for Not Comparing with Shadow Model-Based Attacks

As mentioned in our contributions, IMIA does not need any training data or train shadow models. But regarding the state-of-the-art MIA methods, whether it is LiRA[6], RMIA[34], Oslo[21] or YOQO[31], all of these methods require training additional models. In LiRA[6], RMIA[34], and YOQO[31], shadow models are trained, and in Oslo[21], source and validation models are trained. Although some methods[21] can conduct attacks using a small number of shadow models, the fact is that they still train shadow models. Moreover, all methods based on training shadow models require large amounts of data, which is acknowledged in Prashanth et al. [21]. In our method, no additional models need to be trained, and there is no requirement to know training data to perform IMIA. Therefore, to ensure a fair comparison, we also choose these metric-based methods like Softmax and Entropy that do not require training shadow models as our baselines.

Table 6: The comparison results of computational time cost between IMIA and those methods using shadow models. Time reported under the attack settings in our paper refers to the generation of a single adversarial example per sample. Time reported for RMIA/LiRA refers to the training time of a single shadow model. We compute different adversarial strategies under different model structures and datasets.

| Attack Setting | Attack Method | Model-Data | Time (ms/sample) |
|---|---|---|---|
| Score-based attack | SimBA | ResNet-CIFAR10 | 709.66 |
| Decision-based (label-only) attack | HSJA | DenseNet-STL10 | 3767.1 |
| White-box attack | PGD | VGG-CIFAR100 | 33.41 |
| RMIA / LiRA | Shadow model | ResNet-CIFAR10 | 30min |

## C  The computational time cost of IMIA

In our paper, our IMIA method determines membership by counting the number of iterations required to generate adversarial samples. This process is carried out on a per-sample basis and, since generating adversarial samples in a black-box setting does not require storing gradient information—only forward propagation—the time spent on a single sample is relatively small. We list the specific time cost of computation for each adversarial strategy in Table 6. For example, in the white-box setting, achieving an adversarial sample only needs 33.4 milliseconds under the VGG model trained with CIFAR100. By contrast, methods that rely on training shadow models, even though the models can be reused, typically take at least half an hour to train a single shadow model on an RTX 4090 GPU. When the sample size is large, IMIA may take longer. However, our attack does not require a substantial amount of training data to be effective.