

---

# White-box vs Black-box: Bayes Optimal Strategies for Membership Inference

---

Alexandre Sablayrolles<sup>1,2</sup> Matthijs Douze<sup>2</sup> Yann Ollivier<sup>2</sup> Cordelia Schmid<sup>1</sup> Hervé Jégou<sup>2</sup>

## Abstract

Membership inference determines, given a sample and trained parameters of a machine learning model, whether the sample was part of the training set. In this paper, we derive the optimal strategy for membership inference with a few assumptions on the distribution of the parameters. We show that optimal attacks only depend on the loss function, and thus black-box attacks are as good as white-box attacks. As the optimal strategy is not tractable, we provide approximations of it leading to several inference methods, and show that existing membership inference methods are coarser approximations of this optimal strategy. Our membership attacks outperform the state of the art in various settings, ranging from a simple logistic regression to more complex architectures and datasets, such as ResNet-101 and Imagenet.

## 1. Introduction

Ateniese et al. (2015) state that “*it is unsafe to release trained classifiers since valuable information about the training set can be extracted from them*”. The problem that we address in this paper, *i.e.*, to determine whether a sample has been used to train a given model, is related to the privacy implications of machine learning systems. They were first discussed in the context of support vector machines (Rubinstein et al., 2009; Biggio et al., 2014). The problem of “unintended memorization” (Carlini et al., 2018) appears in most applications of machine learning, such as natural language processing systems (Carlini et al., 2018) or image classification (Yeom et al., 2018).

More specifically, we consider the problem of membership inference, *i.e.*, we aim at determining if a specific image was used to train a model, given only the (image, label) pair and the model parameters. This question is important to protect

both the privacy and intellectual property associated with data. For neural networks, the privacy issue was recently considered by Yeom et al. (2018) for the MNIST and CIFAR datasets. The authors evidence the close relationship between overfitting and privacy of training images. This is reminiscent of prior membership inference attacks, which employ the output of the classifier associated with a particular sample to determine whether it was used during training or not (Shokri et al., 2017).

At this stage, it is worth defining the different levels of information to which the “attacker”, *i.e.*, the membership inference system, has access to. We assume that the attacker knows the data distribution and the specifications of the model (training procedure, architecture of the network, etc), even though they are not necessarily required for all methods. We refer to the *white-box* setting as the case where the attacker knows all the network parameters. On a side note, the setup commonly adopted in differential privacy (Dwork et al., 2006) corresponds to the white-box setting, where the attacker additionally knows all the training samples except the one to be tested.

The *black-box* setting is when these parameters are unknown. For classification models, the attacker has only access to the output for a given input, in one of the following forms:

- (i) the classifier decision;
- (ii) the loss of the correct label;
- (iii) the full response for all classes.

Prior works on membership inference commonly assume (i) or (iii). Our paper focuses on the black-box case (ii), in which we know the loss incurred by the correct label. The state of the art in this setting are the shadow models proposed by Shokri et al. (2017).

In our work, we use a probabilistic framework to derive a formal analysis of the optimal attack. This framework encompasses both Bayesian learning, and noisy training, where the noise is injected (Welling & Teh, 2011) or comes from the stochasticity of SGD. Under mild assumptions on the distribution of the parameters, we derive the optimal membership inference strategy. This strategy only depends on the classifier through evaluation of the loss, thereby showing that black-box attacks will perform as well as white-box attacks in this optimal asymptotic setting. This result may explain why, to the best of our knowledge, the literature

---

<sup>1</sup>University Grenoble Alpes, Inria, CNRS, Grenoble INP, LJK

<sup>2</sup>Facebook AI Research. Correspondence to: Alexandre Sablayrolles <asablayrolles@fb.com>.

does not report white-box attacks outperforming the state-of-the-art black-box-(ii) attacks.

The aforementioned optimal strategy is not tractable, therefore we introduce approximations to derive an explicit method for membership inference. As a byproduct of this derivation, we show that state-of-the-art approaches (Shokri et al., 2017; Yeom et al., 2018) are coarser approximations of the optimal strategy. One of the approximation drastically simplifies the membership inference procedure by simply relying on the loss and a calibration term. We employ this strategy to the more complex case of neural networks, and show that it outperforms all approaches we are aware of.

In summary, our main contributions are as follows:

- We show, under a few assumptions on training, that the optimal inference only depends on the loss function, and not on the parameters of the classifier. In other terms, white-box attacks don't provide any additional information and result in the same optimal strategy.
- We employ different approximations to derive three explicit membership attack strategies. We show that state-of-the-art methods constitute other approximations. Simple simulations show the superiority of our approach on a simple regression problem.
- We apply a simplified, tractable, strategy to infer the membership of images to the train set in the case of the public image classification benchmarks CIFAR and Imagenet. It outperforms the state of the art for membership inference, namely the shadow models.

The paper is organized as follows. Section 2 reviews related work. Section 3 introduces our probabilistic formulation and derives our main theoretical result. This section also discusses the connection between membership inference and differential privacy. Section 4 considers approximations for practical use-cases, which allow us to derive inference strategies in closed-form, some of which are connected with existing methods from the literature. Section 5 summarizes the practical algorithms derived from our analysis. Finally, Section 6 considers the more difficult case of membership inference for real-life neural networks and datasets.

## 2. Related work

**Leakage of information from the training set.** Our work is related to the topics of overfitting and memorization capabilities of classifiers. Determining what neural networks actually memorize from their training set is not trivial. A few recent works (Zhang et al., 2017; Yeom et al., 2018) evaluate how a network can fit random labels. Zhang et al. (2017) replace true labels by random labels and show that

popular neural nets can perfectly fit them in simple cases, such as small datasets (CIFAR10) or Imagenet without data augmentation. Krueger et al. (2017) extend their analysis and argue in particular that the effective capacity of neural nets depends on the dataset considered. In a privacy context, Yeom et al. (2018) exploit this memorizing property to watermark networks. As a side note, random labeling and data augmentation have been used for the purpose of training a network without any annotated data (Dosovitskiy et al., 2014; Bojanowski & Joulin, 2017).

In the context of differential privacy (Dwork et al., 2006), recent works (Wang et al., 2016; Bassily et al., 2016) suggest that guaranteeing privacy requires learning systems to generalize well, *i.e.*, to not overfit. Wang et al. (2015) show that Bayesian posterior sampling offers differential privacy guarantees. Abadi et al. (2016) introduce noisy SGD to learn deep models with differential privacy.

**Membership Inference.** A few recent works (Hayes et al., 2017; Shokri et al., 2017; Long et al., 2018) have addressed membership inference. Yeom et al. (2018) propose a series of membership attacks and derive their performance. Long et al. (2018) observe that some training images are more vulnerable than others and propose a strategy to identify them. Hayes et al. (2017) analyze privacy issues arising in generative models. Dwork et al. (2015) and Sankaranarayanan et al. (2009) provide optimal strategies for membership inference in genomics data.

**Shadow models.** Shadow models were introduced by Shokri et al. (2017) in the context of black-box attacks. In this setup, an attacker has black-box-(iii) access (full response for all classes) to a model trained on a private dataset, and to a public dataset that follows the same distribution as the private dataset. The attacker wishes to perform membership inference using black-box outputs of the private model. For this, the attacker simulates models by training *shadow models* on known splits from the public set. On this simulated models, the attacker can analyze the output patterns corresponding to samples from the training set and from a held-out set. Shokri et al. (2017) propose to train an attack model that learns to predict, given an output pattern, whether it corresponds to a training or held-out sample. If the attack model simply predicts “training” when the output activations fire on the correct class, this strategy is equivalent to Yeom et al. (2018)’s adversary. Salem et al. (2019) further show that shadow models work under weaker assumptions than those of Shokri et al. (2017).

## 3. Membership inference model

In this section, we derive the Bayes optimal performance for membership inference (Theorem 1). We then make the connection with differential privacy and propose looser guarantees that prevent membership inference.

### 3.1. Posterior distribution of parameters

Let  $\mathcal{D}$  be a data distribution, from which we sample  $n \in \mathbb{N}$  points  $z_1, z_2, \dots, z_n \underset{\text{i.i.d.}}{\sim} \mathcal{D}$ . A machine learning algorithm produces parameters  $\theta$  that incur a low loss  $\sum_{i=1}^n \ell(\theta, z_i)$ . Typically in the case of classification,  $z = (x, y)$  where  $x$  is an input sample and  $y$  a class label, and the loss function  $\ell$  is high on samples  $(x, y')$  for which  $y' \neq y$ .

We assume that the machine learning algorithm has some randomness, and we model it with a posterior distribution over parameters  $\theta | z_1, \dots, z_n$ . The randomness in  $\theta$  either comes from the training procedure (e.g., Bayesian posterior sampling), or arises naturally, as is the case with Stochastic Gradient methods.

In general, we assume that the posterior distribution follows:

$$\mathbb{P}(\theta | z_1, \dots, z_n) \propto e^{-\frac{1}{T} \sum_{i=1}^n \ell(\theta, z_i)}, \quad (1)$$

where  $T$  is a temperature parameter, intuitively controlling the stochasticity of  $\theta$ .  $T = 1$  corresponds to the case of the Bayesian posterior,  $T \rightarrow 0$  the case of MAP (Maximum A Posteriori) inference and a small  $T$  to the case of averaged SGD (Polyak & Juditsky, 1992). Note that we do not include a prior on  $\theta$ : we assume that the prior is either uniform on a bounded  $\theta$ , or that it has been factored in the loss term  $\ell$ .

### 3.2. Membership inference

Given  $\theta$  produced by such a machine learning algorithm, membership inference asks the following question: What information does  $\theta$  contain about its training set  $z_1, \dots, z_n$ ?

Formally, we assume that binary membership variables  $m_1, m_2, \dots, m_n$  are drawn independently, with probability  $\lambda = \mathbb{P}(m_i = 1)$ . The samples for which  $m_i = 0$  are the test set, while the samples for which  $m_i = 1$  are the training set. Equation (1) becomes

$$\mathbb{P}(\theta | z_1, \dots, z_n, m_1, \dots, m_n) \propto e^{-\frac{1}{T} \sum_{i=1}^n m_i \ell(\theta, z_i)}, \quad (2)$$

Taking the case of  $z_1$  without loss of generality, membership inference determines, given parameters  $\theta$  and sample  $z_1$ , whether  $m_1 = 1$  or  $m_1 = 0$ .

**Definition 1** (Membership inference). *Inferring the membership of sample  $z_1$  to the training set amounts to computing:*

$$\mathcal{M}(\theta, z_1) := \mathbb{P}(m_1 = 1 | \theta, z_1). \quad (3)$$

**Notation.** We denote by  $\sigma$  the sigmoid function  $\sigma(u) = (1 + e^{-u})^{-1}$ . We collect the knowledge about the other samples and their memberships into the set  $\mathcal{T} = \{z_2, \dots, z_n, m_2, \dots, m_n\}$ .

### 3.3. Optimal membership inference

In Theorem 1, we derive the explicit formula for  $\mathcal{M}(\theta, z_1)$ .

**Theorem 1.** *Given a parameter  $\theta$  and a sample  $z_1$ , the optimal membership inference is given by:*

$$\mathcal{M}(\theta, z_1) = \mathbb{E}_{\mathcal{T}} \left[ \sigma \left( \log \left( \frac{\mathbb{P}(\theta | m_1 = 1, z_1, \mathcal{T})}{\mathbb{P}(\theta | m_1 = 0, z_1, \mathcal{T})} \right) + t_\lambda \right) \right], \quad (4)$$

with  $t_\lambda = \log \left( \frac{\lambda}{1-\lambda} \right)$ .

*Proof.* By the law of total expectation, we have:

$$\mathcal{M}(\theta, z_1) = \mathbb{P}(m_1 = 1 | \theta, z_1) \quad (5)$$

$$= \mathbb{E}_{\mathcal{T}} [\mathbb{P}(m_1 = 1 | \theta, z_1, \mathcal{T})]. \quad (6)$$

Applying Bayes' formula:

$$\mathbb{P}(m_1 = 1 | \theta, z_1, \mathcal{T}) = \frac{\mathbb{P}(\theta | m_1 = 1, z_1, \mathcal{T}) \mathbb{P}(m_1 = 1)}{\mathbb{P}(\theta | z_1, \mathcal{T})}, \quad (7)$$

$$= \frac{\alpha}{\alpha + \beta} = \sigma \left( \log \left( \frac{\alpha}{\beta} \right) \right) \quad (8)$$

where:

$$\alpha := \mathbb{P}(\theta | m_1 = 1, z_1, \mathcal{T}) \mathbb{P}(m_1 = 1) \quad (9)$$

$$\beta := \mathbb{P}(\theta | m_1 = 0, z_1, \mathcal{T}) \mathbb{P}(m_1 = 0) \quad (10)$$

Given that  $\mathbb{P}(m_1 = 1) = \lambda$ ,

$$\log \left( \frac{\alpha}{\beta} \right) = \log \left( \frac{\mathbb{P}(\theta | m_1 = 1, z_1, \mathcal{T})}{\mathbb{P}(\theta | m_1 = 0, z_1, \mathcal{T})} \right) + \log \left( \frac{\lambda}{1-\lambda} \right), \quad (11)$$

which gives the expression for  $\mathcal{M}(\theta, z_1)$ .  $\square$

Note that Theorem 1 only relies on the fact that  $\theta$  given  $\{z_1, \dots, z_n, m_1, \dots, m_n\}$  is a random variable, but it does not make any assumption on the form of the distribution. In particular the loss  $\ell$  does not appear in the expression.

Theorem 2 uses the assumption in Equation (2) to further explicit  $\mathcal{M}(\theta, z_1)$ ; we give its formal expression below, prove it, and analyze the expression qualitatively. Let us first define the posterior over the parameters given samples  $z_2, \dots, z_n$  and memberships  $m_2, \dots, m_n$ :

$$p_{\mathcal{T}}(\theta) := \frac{e^{-\frac{1}{T} \sum_{i=2}^n m_i \ell(\theta, z_i)}}{\int_t e^{-\frac{1}{T} \sum_{i=2}^n m_i \ell(t, z_i)} dt}. \quad (12)$$

**Theorem 2.** *Given a parameter  $\theta$  and a sample  $z_1$ , the optimal membership inference is given by:*

$$\mathcal{M}(\theta, z_1) = \mathbb{E}_{\mathcal{T}} [\sigma (s(z_1, \theta, p_{\mathcal{T}}) + t_\lambda)] \quad (13)$$

where we define the following score:

$$\tau_p(z_1) := -T \log \left( \int_t e^{-\frac{1}{T} \ell(t, z_1)} p(t) dt \right) \quad (14)$$

$$s(z_1, \theta, p) := \frac{1}{T} (\tau_p(z_1) - \ell(\theta, z_1)). \quad (15)$$

*Proof.* Singling out  $m_1$  in Equation (2) yields the following expressions for  $\alpha$  and  $\beta$ :

$$\alpha = \lambda \frac{e^{-\frac{1}{T} \ell(\theta, z_1)} e^{-\frac{1}{T} \sum_{i=2}^n m_i \ell(\theta, z_i)}}{\int_t e^{-\frac{1}{T} \ell(t, z_1)} e^{-\frac{1}{T} \sum_{i=2}^n m_i \ell(t, z_i)} dt} \quad (16)$$

$$= \lambda \frac{e^{-\frac{1}{T} \ell(\theta, z_1)} p_{\mathcal{T}}(\theta)}{\int_t e^{-\frac{1}{T} \ell(t, z_1)} p_{\mathcal{T}}(t) dt}, \quad (17)$$

and

$$\beta = (1 - \lambda) \frac{e^{-\frac{1}{T} \sum_{i=2}^n m_i \ell(\theta, z_i)}}{\int_t e^{-\frac{1}{T} \sum_{i=2}^n m_i \ell(t, z_i)} dt} = (1 - \lambda) p_{\mathcal{T}}(\theta). \quad (18)$$

Thus,

$$\begin{aligned} \log \left( \frac{\alpha}{\beta} \right) &= -\frac{\ell(\theta, z_1)}{T} - \log \left( \int_t e^{-\frac{1}{T} \ell(t, z_1)} p_{\mathcal{T}}(t) dt \right) + t_{\lambda} \\ &= s(z_1, \theta, p_{\mathcal{T}}) + t_{\lambda}. \end{aligned} \quad (19)$$

Then, Equation (8) yields the expected result.  $\square$

The first observation in Theorem 2 is that  $\mathcal{M}(\theta, z_1)$  does not depend on the parameters  $\theta$  beyond the evaluation of the loss  $\ell(\theta, \cdot)$ : this strategy does not require, for instance, internal parameters of the model that a white-box attack could provide. This means that if we can compute  $\tau_p$  or approximate it well enough, then the optimal membership inference depends only on the loss. In other terms, asymptotically, **the white-box setting does not provide any benefit compared to black-box membership inference.**

Let us analyze qualitatively the terms in the expression: Since  $\mathcal{T}$  is a training set,  $p_{\mathcal{T}}$  corresponds to a posterior over this training set, *i.e.*, a typical distribution of trained parameters.  $\tau_p(z_1)$  is the softmin of loss terms  $\ell(\cdot, z_1)$  over these typical models, and corresponds therefore to the *typical loss of sample*  $z_1$  under models that have not seen  $z_1$ .

The quantity  $\tau_p(z_1)$  can be seen as a threshold, to which the loss  $\ell(\theta, z_1)$  is compared. Around this threshold, when  $\ell(\theta, z_1) \approx \tau_p(z_1)$ , then  $s \approx 0$ : since  $\sigma(t_{\lambda}) = \lambda$ , the membership posterior probability  $\mathcal{M}(\theta, z_1)$  is equal to  $\lambda$ , and thus we have no information on  $m_1$  beyond prior knowledge. As the loss  $\ell(\cdot, z_1)$  gets lower than this threshold,  $s$  becomes positive. Since  $\sigma$  is non decreasing, when  $s(z_1, \theta, p_{\mathcal{T}}) > 0$ ,  $\mathcal{M}(\theta, z_1) > \lambda$  and thus we gain non-trivial membership information on  $z_1$ .

Another consequence of Theorem 2 is that a higher temperature  $T$  decreases  $s$ , and thus decreases  $\mathbb{P}(m_1 = 1 \mid \theta, z_1)$ : it corresponds to the intuition that more randomness in  $\theta$  protects the privacy of training data.

### 3.4. Differential privacy and guarantees

In this subsection we make the link with differential privacy.

Differential privacy (Dwork et al., 2006) is a framework that allows to learn model parameters  $\theta$  while maintaining the confidentiality of data. It ensures that even if a malicious attacker knows parameters  $\theta$  and samples  $z_i$ ,  $i \geq 2$ , for which  $m_i = 1$ , the privacy of  $z_1$  is not compromised.

**Definition 2** ( $\epsilon$ -differential privacy). *A machine learning algorithm is  $\epsilon$ -differentially private if, for any choice of  $z_1$  and  $\mathcal{T}$ ,*

$$\log \left( \frac{\mathbb{P}(\theta \mid m_1 = 1, z_1, \mathcal{T})}{\mathbb{P}(\theta \mid m_1 = 0, z_1, \mathcal{T})} \right) < \epsilon. \quad (20)$$

Note that this definition is slightly different from the one of Dwork et al. (2006) in that we consider the removal of  $z_1$  rather than its substitution with  $z'$ . Additionally we consider probability densities instead of probabilities of sets, without loss of generality.

**Property 1** ( $\epsilon$ -differential privacy). *If the training is  $\epsilon$ -differentially private, then:*

$$\mathbb{P}(m_1 = 1 \mid \theta, z_1) \leq \lambda + \frac{\epsilon}{4}. \quad (21)$$

*Proof.* Combining Equation (20) and the fact that  $\sigma(u) \leq \sigma(v) + \max(u - v, 0)/4$  (Appendix A.3), we have:

$$\begin{aligned} \sigma \left( \log \left( \frac{\mathbb{P}(\theta \mid m_1 = 1, z_1, \mathcal{T})}{\mathbb{P}(\theta \mid m_1 = 0, z_1, \mathcal{T})} \right) + t_{\lambda} \right) &\leq \sigma(t_{\lambda}) + \frac{\epsilon}{4} \\ &= \lambda + \frac{\epsilon}{4}. \end{aligned} \quad (22)$$

Combining this expression with Theorem 1 yields the result.  $\square$

Note that this bound gives a tangible sense of  $\epsilon$ . In general, decreasing  $\epsilon$  increases privacy, but there is no consensus over “good” values of  $\epsilon$ ; this bound indicates for instance that  $\epsilon = 0.01$  would be sufficient for membership privacy.

$\epsilon$ -differential privacy gives strong membership inference guarantees, at the expense of a constrained training procedure resulting generally in a loss of accuracy (Abadi et al., 2016). However, if we assume that the attacker knows the  $z_i, i \geq 2$  for which  $m_i = 1$ ,  $\epsilon$ -differential privacy is required to protect the privacy of  $z_1$ . Depending on the information we have on  $z_i, i \geq 2$ , there is a continuum between differential privacy (all  $z_i$ 's are known) and membership inference (only prior knowledge on  $z_i$ ). In the case of membership inference, it suffices to have the following guarantee:

**Definition 3** ( $(\epsilon, \delta)$  membership privacy). *The training is  $(\epsilon, \delta)$ -membership private for some  $\epsilon > 0, \delta > 0$  if with probability  $1 - \delta$  over the choice of  $\mathcal{T}$ :*

$$\int_t \ell(t, z_1) p_{\mathcal{T}}(t) dt - \ell(\theta, z_1) \leq \epsilon. \quad (23)$$

**Property 2.** *If the training is  $(\epsilon, \delta)$ -membership private, then:*

$$\mathbb{P}(m_1 = 1 \mid \theta, z_1) \leq \lambda + \frac{\epsilon}{4T} + \delta. \quad (24)$$

*Proof.* Jensen’s inequality states that for any distribution  $p$  and any function  $f$ :

$$\int_t f(t) p(t) dt \leq \log \left( \int_t e^{f(t)} p(t) dt \right), \quad (25)$$

hence the score  $s$  from Equation (15) verifies:

$$s(z_1, \theta, p) \leq \frac{1}{T} \left( \int_t \ell(t, z_1) p(t) dt - \ell(\theta, z_1) \right). \quad (26)$$

Thus, distinguishing the cases  $\delta$  and  $1 - \delta$  in the expectation in Equation (13),

$$\mathbb{P}(m_1 = 1 \mid \theta, z_1) \leq \delta + (1 - \delta) \left( \lambda + \frac{\epsilon}{4T} \right) \quad (27)$$

$$\leq \lambda + \frac{\epsilon}{4T} + \delta, \quad (28)$$

which gives the desired bound.  $\square$

Membership privacy provides a *post-hoc* guarantee on  $\theta, z_1$ . Guarantees in the form of Equation (23) can be obtained by PAC (Probably Approximately Correct) bounds.

## 4. Approximations for membership inference

Estimating the probability of Equation (15) mainly requires to compute the term  $\tau_p$ . Since its expression is intractable, we use approximations to derive concrete membership attacks (MA). We now detail these approximations, referred to as MAST (MA Sample Threshold), MALT (MA Loss Threshold) and MATT (MA Taylor Threshold).

### 4.1. MAST: Approximation of $\tau(z_1)$

We first make the mean-field assumption that  $p_{\mathcal{T}}(t)$  does not depend on  $\mathcal{T}$  (we note it  $p$ ), and define

$$\tau(z_1) := \log \left( \int_t e^{-\frac{1}{T} \ell(t, z_1)} p(t) dt \right). \quad (29)$$

The quantity  $\tau(\cdot)$  is a “calibrating” term that reflects the difficulty of a sample. Intuitively, a low  $\tau(z_1)$  means that the sample  $z_1$  is easy to predict, and thus a low value of  $\ell(\theta, z_1)$  does not necessarily indicate that  $z_1$  belongs to the train set. Thus, Theorem 2 gives the optimal attack model:

$$s_{\text{MAST}}(\theta, z_1) = -\ell(\theta, z_1) + \tau(z_1). \quad (30)$$

### 4.2. MALT: Constant $\tau$

If we further assume that  $\tau(\cdot)$  is constant, the optimal strategy reduces to predicting that  $z_1$  comes from the training set if its loss  $\ell(\theta, z_1)$  is lower than this threshold  $\tau$ , and from the test set otherwise:

$$s_{\text{MALT}}(\theta, z_1) = -\ell(\theta, z_1) + \tau. \quad (31)$$

A similar strategy is proposed by Yeom et al. (2018) for Gaussian models. Carlini et al. (2018) estimate a secret token in text datasets by comparing probabilities of the sentence “My SSN is X” with various values of X. Surprisingly, to the best of our knowledge, it has not been proposed in the literature to estimate the threshold  $\tau$  on public data, and to apply it for membership inference. As we show in Section 6, this simple strategy yields better results than shadow models. Their attack models take as input the softmax activation  $\phi_{\theta}(x)$  and the class label  $y$ , and predict whether the sample comes from the training or test set. For classification models,  $\ell(\theta, (x, y)) = -\log(\phi_{\theta}(x)_y)$ . Hence the optimal attack performs:

$$s_{\text{MALT}}(\theta, (x, y)) = \log(\phi_{\theta}(x)_y) + \tau. \quad (32)$$

In Shokri et al. (2017), we argue that the attack model essentially performs such an estimation, albeit in a non-explicit way. In particular, we believe that the gap between Shokri et al. (2017)’s method and ours is due to instabilities in the estimation of  $\tau$  and the numerical computation of the log, as the model is given only  $\phi_{\theta}(x)$ . As a side note, the expectation term in  $\mathcal{T} = z_2, \dots, z_n, m_2, \dots, m_n$  is very similar in spirit to the shadow models, and they can be viewed as a Monte-Carlo estimation of this quantity.

**An experiment with Gaussian data.** We illustrate the difference between a MALT (global  $\tau$ ) and MAST (per-sample  $\tau(\cdot)$ ) on a simple toy example. Let’s assume we estimate the mean  $\mu$  of Gaussian data with unit variance.

We sample  $n$  values  $z_1, \dots, z_n$  from  $\mathcal{D} = \mathcal{N}(\mu, I)$ . The estimate of the mean is  $\theta = \frac{1}{n'} \sum_{i=1}^n m_i z_i$  where  $n' = |\{i \mid m_i = 1\}|$ . We have (see Appendix A.2 for derivations):

$$\ell(\theta, z_i) := \frac{1}{2} \|z_i - \theta\|^2 \quad (33)$$

$$\tau(z_i) = \frac{n'}{2(n'+1)} \|z_i - \mu\|^2 \quad (34)$$

$$\tau = \frac{n'}{2(n'+1)} \mathbb{E} \|z - \mu\|^2 = \frac{n'}{2(n'+1)} d. \quad (35)$$

The expression of  $\tau(z_i)$  shows that the “difficulty” of sample  $z_i$  is its distance to  $\mu$ , *i.e.*, how untypical this sample is.

Figure 1 shows the results with a global  $\tau$  or a per-sample  $\tau$ : the per-sample  $\tau$  better separates the two distributions, leading to an increased membership inference accuracy.

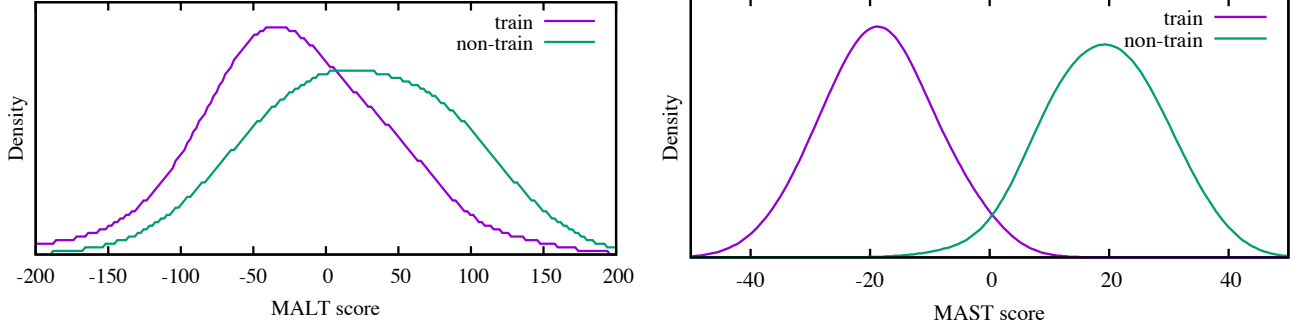


Figure 1. Comparison of MALT and MALT for membership inference on the mean estimator for Gaussian data ( $n = 100$  samples in 2000 dimensions). Distribution of scores  $s$  used to distinguish between samples seen or not at training. *MALT*: a single threshold is used for all the dataset; *MAST*: each sample gets assigned a different threshold. *MAST* better separates training and non-training samples.

**MATT: Estimation with Taylor expansion** We assume that the posterior induced by the loss  $\mathcal{L}(\theta) = \sum_{i=1}^n m_i \ell(\theta, z_i)$  is a Gaussian centered on  $\theta^*$ , the minimum of the loss, with covariance  $C$ . This corresponds to the Laplace approximation of the posterior distribution. The inverse covariance matrix  $C^{-1}$  is asymptotically  $n$  times the Fisher matrix (van der Vaart, 1998), which itself is the Hessian of the loss (Kullback, 1997):

$$\mathbb{P}(\theta | z_1, \mathcal{T}) = \frac{1}{\sqrt{\det(2\pi H^{-1})}} e^{-\frac{1}{2}(\theta - \theta^*)^T H (\theta - \theta^*)}. \quad (36)$$

We denote by  $\theta_0^*$  (resp.  $\theta_1^*$ ) the mode of the Gaussian corresponding to  $\{z_2, \dots, z_n\}$  (resp.  $\{z_1, \dots, z_n\}$ ), and  $H_0$  (resp.  $H_1$ ) the corresponding Hessian matrix. We assume that  $H$  is not impacted by removing  $z_1$  from the training set, and thus  $H := H_0 \approx H_1$  (cf. appendix A.4 for a more precise justification). The log-ratio is therefore

$$\begin{aligned} \log \left( \frac{\mathbb{P}(\theta | m_1 = 1, z_1, \mathcal{T})}{\mathbb{P}(\theta | m_1 = 0, z_1, \mathcal{T})} \right) & \quad (37) \\ &= -\frac{1}{2}(\theta - \theta_1^*)^T H (\theta - \theta_1^*) + \frac{1}{2}(\theta - \theta_0^*)^T H (\theta - \theta_0^*) \\ &= (\theta - \theta_0^*)^T H (\theta_1^* - \theta_0^*) - \frac{1}{2}(\theta_1^* - \theta_0^*)^T H (\theta_1^* - \theta_0^*). \end{aligned}$$

The difference  $\theta_1^* - \theta_0^*$  can be estimated using a Taylor expansion of the loss gradient around  $\theta_0^*$  (see e.g. Koh & Liang (2017)):

$$\theta_1^* - \theta_0^* \approx -H^{-1} \nabla_{\theta} \ell(\theta_0^*, z_1) \quad (38)$$

Combining this with Equation (37) leads to

$$-(\theta - \theta_0^*)^T \nabla_{\theta} \ell(\theta_0^*, z_1) - \frac{1}{2} \nabla_{\theta} \ell(\theta_0^*, z_1)^T H^{-1} \nabla_{\theta} \ell(\theta_0^*, z_1). \quad (39)$$

We study the asymptotic behavior of this expression when  $n \rightarrow \infty$ . On the left-hand side, the parameters  $\theta$  and  $\theta_0^*$  are estimates of the optimal  $\theta^*$ , and under mild conditions, the

error of the estimated parameters is of order  $1/\sqrt{n}$ . Therefore the difference  $\theta - \theta_0^*$  is of order  $1/\sqrt{n}$ . On the right-hand side, the matrix  $H$  is the summation of  $n$  sample-wise Hessian matrices. Therefore, asymptotically, the right-hand side shrinks at a rate  $1/n$ , which is negligible compared to the other, which shrinks at  $1/\sqrt{n}$ . In addition to the asymptotic reasoning, we verified this approximation experimentally. Thus, we approximate Equation (39) to give the following score:

$$s_{\text{MATT}}(\theta, z_1) = -(\theta - \theta_0^*)^T \nabla_{\theta} \ell(\theta_0^*, z_1). \quad (40)$$

Equation (40) has an intuitive interpretation: parameters  $\theta$  were trained using  $z_1$  if their difference with a set of parameters trained without  $z_1$  (i.e.  $\theta_0^*$ ) is aligned with the direction of the update  $-\nabla_{\theta} \ell(\theta_0^*, z_1)$ .

## 5. Membership inference algorithms

In this section, we detail how the approximations of  $s(\theta, z_1, p)$  are employed to perform membership inference.

We assume that a machine learning model has been trained, yielding parameters  $\theta$ . We assume also that similar models can be re-trained with different training sets. Given a sample  $z_1$ , we want to decide whether  $z_1$  belongs to the training set.

### 5.1. The 0-1 baseline

We consider as a baseline the “0-1” heuristic, which predicts that  $z_1$  comes from the training set if the class is predicted correctly, and from the test set if the class is predicted incorrectly. We note  $p_{\text{train}}$  (resp.  $p_{\text{test}}$ ) the classification accuracy on the training (resp. held-out) set. The accuracy of the heuristic is (see Appendix A.1 for derivations):

$$p_{\text{bayes}} = \lambda p_{\text{train}} + (1 - \lambda)(1 - p_{\text{test}}). \quad (41)$$

For example when  $\lambda = 1/2$ , since  $p_{\text{train}} \geq p_{\text{test}}$  this heuristic

is better than random guessing (accuracy  $1/2$ ) and the improvement is proportional to the overfitting gap  $p_{\text{train}} - p_{\text{test}}$ .

## 5.2. Making hard decisions from scores

Variants of our method provide different estimates of  $s(\theta, z_1, p)$ . Theorem 2 shows that this score has to be passed through a sigmoid function, but since it is an increasing function, the threshold can be chosen directly on these scores. Estimation of this threshold has to be conducted on simulated sets, for which membership information is known. We observed that there is almost no difference between choosing the threshold on the set to be tested and cross-validating it. This is expected, as a one-dimensional parameter the threshold is not prone to overfitting.

## 5.3. Membership algorithms

**MALT: Threshold on the loss.** Since  $\tau$  in Equation (31) is constant, and using the invariance to increasing functions, we need only to use loss value for the sample,  $\ell(\theta, z_1)$ .

**MAST: Estimating  $\tau(z_1)$ .** To estimate  $\tau(z_1)$  in Equation (29), we train several models with different subsamples of the training set. This yields a set of per-sample losses for  $z_1$  that are averaged into an estimate of  $\tau(z_1)$ .

**MATT: the Taylor approximation.** We run the training on a separate set to obtain  $\theta_0^*$ . Then we take a gradient step over the loss to estimate the approximation in Equation (40). Note that this strategy is not compatible with neural networks because the assumption that parameters lie around a unique global minimum does not hold. In addition, parameters from two different networks  $\theta$  and  $\theta_0^*$  cannot be compared directly as neural networks that express the same function can have very different parameters (*e.g.* because channels can be permuted arbitrarily).

## 6. Experiments

In this section we evaluate the membership inference methods on machine-learning tasks of increasing complexity.

### 6.1. Evaluation

We evaluate three metrics: the accuracy of the attack, and the mean average precision when detecting either from train ( $\text{mAP}_{\text{train}}$ ) or test ( $\text{mAP}_{\text{test}}$ ) images. For the mean average precision, the scores need not to be thresholded, the metric is invariant to mapping by any increasing function.

### 6.2. Logistic regression

CIFAR-10 is a dataset of  $32 \times 32$  pixel images grouped in 10 classes. In this subsection, we consider two of the classes (*truck* and *boat*) and vary the number of training images

Table 1. Accuracy (top) and mAP (bottom) of membership inference on the 2-class logistic regression with simple CNN features, for different types of attacks. Note that 0-1 corresponds to the baseline (Yeom et al., 2018). We do not report Shokri et al. (2017) since Table 2 shows MALT performs better. Results are averaged over 100 different random seeds.

| n    | Model accuracy |            | Attack accuracy |      |             |
|------|----------------|------------|-----------------|------|-------------|
|      | train          | validation | 0 – 1           | MALT | MATT        |
| 400  | 97.9           | 93.8       | 52.1            | 54.4 | <b>57.0</b> |
| 1000 | 97.3           | 94.5       | 51.4            | 52.6 | <b>54.5</b> |
| 2000 | 96.8           | 95.2       | 50.8            | 51.7 | <b>53.0</b> |
| 4000 | 97.7           | 95.6       | 51.0            | 51.4 | <b>52.1</b> |
| 6000 | 97.5           | 96.0       | 50.7            | 51.0 | <b>51.8</b> |

| n    | $\text{mAP}_{\text{test}}$ |      | $\text{mAP}_{\text{train}}$ |      |
|------|----------------------------|------|-----------------------------|------|
|      | MALT                       | MATT | MALT                        | MATT |
| 400  | 55.8                       | 60.1 | 51.9                        | 57.1 |
| 1000 | 53.2                       | 56.6 | 50.5                        | 54.8 |
| 2000 | 51.8                       | 54.4 | 50.4                        | 53.4 |
| 4000 | 51.9                       | 53.7 | 50.1                        | 52.6 |
| 6000 | 51.4                       | 53.0 | 50.2                        | 52.2 |

from  $n = 400$  to 6,000.

We train a logistic regression to separate the two classes. The logistic regression takes as input features extracted from a pretrained Resnet18 on CIFAR-100 (a disjoint dataset). The regularization parameter  $C$  of the logistic regression is cross-validated on held-out data.

We assume that  $\lambda = 1/2$  ( $n/2$  training images and  $n/2$  test images). We also reserve  $n/2$  images to estimate  $\theta_0^*$  for the MATT method. In both experiments, we report the peak accuracy obtained for the best threshold (cf. Section 5.2).

Table 1 shows the results of our experiments, in terms of accuracy and mean average precision. In accuracy, the Taylor expansion method MATT outperforms the MALT method, for any number of training instances  $n$ , which itself obtains much better results than the naive 0-1 attack.

Interestingly, it shows a difference between MALT and MATT: both perform similarly in terms of  $\text{mAP}_{\text{test}}$ , but MATT slightly outperforms MALT in  $\text{mAP}_{\text{train}}$ . The main reason for this difference is that the MALT attack is asymmetric: it is relatively easy to predict that elements come from the test set, as they have a high loss, but elements with a low loss can come either from the train set or the test set.

Table 2. Accuracy of membership attacks on the CIFAR-10 classification with a simple neural network. The numbers for the related works are from the respective papers.

| Method                              | Accuracy    |
|-------------------------------------|-------------|
| 0-1 (Yeom et al., 2018)             | 69.4        |
| Shadow models (Shokri et al., 2017) | 73.9        |
| MALT                                | <b>77.1</b> |
| MAST                                | 77.6        |

### 6.3. Small convolutional network

In this section we train a small convolutional network<sup>1</sup> with the same architecture as Shokri et al. (2017) on the CIFAR-10 dataset, using a training set of 15,000 images. Our model is trained for 50 epochs with a learning rate of 0.001. We assume a balanced prior on membership ( $\lambda = 1/2$ ).

We run the MALT and MAST attacks on the classifiers. As stated before, the MATT attack cannot be carried out on convolutional networks. For MAST, the threshold is estimated from 30 shadow models: we train these 30 shadow models on 15,000 images chosen among the train+held-out set (30,000 images). Thus, for each image, we have on average 15 models trained on it and 15 models not trained on it: we estimate the threshold for this image by taking the value  $\tau(z)$  that separates the best the two distributions: this corresponds to a non-parametric estimation of  $\tau(z)$ .

Table 2 shows that our estimations outperform the related works. Note that this setup gives a slight advantage to MAST as the threshold is estimated directly for each sample under investigation, whereas MALT first estimates a threshold, and then applies it to never-seen data. Yet, in contrast with the experiment on Gaussian data, MAST performs only slightly better than MALT. Our interpretation for this is that the images in the training set have a high variability, so it is difficult to obtain a good estimate of  $\tau(z_1)$ . Furthermore, our analysis of the estimated thresholds  $\tau(z_1)$  show that they are very concentrated around a central value  $\tau$ , so their impact when added to the scores is limited.

Therefore, in the following experiment we focus on the MALT attack.

### 6.4. Evaluation on Imagenet

We evaluate a real-world dataset and tackle classification with large neural networks on the Imagenet dataset (Deng et al., 2009; Russakovsky et al., 2015), which contains 1.2 million images partitioned into 1000 categories. We divide Imagenet equally into two splits, use one for training and hold out the rest of the data.

<sup>1</sup> 2 convolutional and 2 linear layers with Tanh non-linearity.

Table 3. Imagenet classification with deep convolutional networks: Accuracy of membership inference attacks of the models.

| Model     | Augmentation       | 0-1  | MALT |
|-----------|--------------------|------|------|
| Resnet101 | None               | 76.3 | 90.4 |
|           | Flip, Crop $\pm 5$ | 69.5 | 77.4 |
|           | Flip, Crop         | 65.4 | 68.0 |
| VGG16     | None               | 77.4 | 90.8 |
|           | Flip, Crop $\pm 5$ | 71.3 | 79.5 |
|           | Flip, Crop         | 63.8 | 64.3 |

We experiment with the popular VGG-16 (Simonyan & Zisserman, 2014) and Resnet-101 (He et al., 2016) architectures. The model is learned in 90 epochs, with an initial learning rate of 0.01, divided by 10 every 30 epochs. Parameter optimization is conducted with SGD with a momentum of 0.9, a weight decay of  $10^{-4}$ , and a batch size of 256.

We conduct the membership inference test by running the 0-1 attack and MALT. An important factor for the success of the attacks is the amount of data augmentation. To assess the effect of data augmentation, we train different networks with varying data augmentation: None, Flip+Crop $\pm 5$ , Flip+Crop (by increasing intensity).

Table 3 shows that data augmentation reduces the gap between the training and the held-out accuracy. This decreases the accuracy of the Bayes attack and the MALT attack. As we can see, without data augmentation, it is possible to guess with high accuracy if a given image was used to train a model (about 90% with our approach, against 77% for existing approaches). Stronger data augmentation reduces the accuracy of the attacks, that still remain above 64%.

## 7. Conclusion

This paper has addressed the problem of membership inference by adopting a probabilistic point of view. This led us to derive the optimal inference strategy. This strategy, while not explicit and therefore not applicable in practice, does not depend on the parameters of the classifier if we have access to the loss. Therefore, a main conclusion of this paper is to show that, asymptotically, white-box inference does not provide more information than an optimized black-box setting.

We then proposed two approximations that lead to three concrete strategies. They outperform competitive strategies for a simple logistic problem, by a large margin for our most sophisticated approach (MATT). Our simplest strategy (MALT) is applied to the more complex problem of membership inference from a deep convolutional network on Imagenet, and significantly outperforms the baseline.



## References

- Martin Abadi, Andy Chu, Ian Goodfellow, Brendan McMahhan, Ilya Mironov, Kunal Talwar, and Li Zhang. Deep learning with differential privacy. In *CCS*, 2016.
- Giuseppe Ateniese, Luigi V Mancini, Angelo Spognardi, Antonio Villani, Domenico Vitali, and Giovanni Felici. Hacking smart machines with smarter ones: How to extract meaningful data from machine learning classifiers. *IJSN*, 2015.
- Raef Bassily, Kobbi Nissim, Adam Smith, Thomas Steinke, Uri Stemmer, and Jonathan Ullman. Algorithmic stability for adaptive data analysis. In *STOC*, 2016.
- Battista Biggio, Iginio Corona, Blaine Nelson, Benjamin I. P. Rubinstein, Davide Maiorca, Giorgio Fumera, Giorgio Giacinto, and Fabio Roli. *Security Evaluation of Support Vector Machines in Adversarial Environments*. Springer International Publishing, 2014.
- Piotr Bojanowski and Armand Joulin. Unsupervised learning by predicting noise. In *ICML*, 2017.
- Nicholas Carlini, Chang Liu, Jernej Kos, Úlfar Erlingsson, and Dawn Song. The secret sharer: Measuring unintended neural network memorization & extracting secrets. *arXiv preprint arXiv:1802.08232*, 2018.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009.
- Alexey Dosovitskiy, Jost Tobias Springenberg, Martin Riedmiller, and Thomas Brox. Discriminative unsupervised feature learning with convolutional neural networks. In *NIPS*, 2014.
- Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity in private data analysis. In *TCC*, 2006.
- Cynthia Dwork, Adam Smith, Thomas Steinke, Jonathan Ullman, and Salil Vadhan. Robust traceability from trace amounts. In *Proceedings of the Symposium on the Foundations of Computer Science*, 2015.
- Jamie Hayes, Luca Melis, George Danezis, and Emiliano De Cristofaro. Logan: evaluating privacy leakage of generative models using generative adversarial networks. *arXiv preprint arXiv:1705.07663*, 2017.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In *ICML*, 2017.
- David Krueger, Nicolas Ballas, Stanislaw Jastrzebski, Devansh Arpit, Maxinder S. Kanwal, Tegan Maharaj, Emmanuel Bengio, Asja Fischer, Aaron Courville, Simon Lacoste-Julien, and Yoshua Bengio. A closer look at memorization in deep networks. In *ICML*, 2017.
- S. Kullback. *Information Theory And Statistics*. Dover Publications, 1997.
- Yunhui Long, Vincent Bindschaedler, Lei Wang, Diyue Bu, Xiaofeng Wang, Haixu Tang, Carl A Gunter, and Kai Chen. Understanding membership inferences on well-generalized learning models. *arXiv preprint arXiv:1802.04889*, 2018.
- B. T. Polyak and A. B. Juditsky. Acceleration of stochastic approximation by averaging. *SIAM J. Control Optim.*, 1992.
- Benjamin IP Rubinstein, Peter L Bartlett, Ling Huang, and Nina Taft. Learning in a large function space: Privacy-preserving mechanisms for SVM learning. *arXiv:0911.5708*, 2009.
- Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *IJCV*, 2015.
- Ahmed Salem, Yang Zhang, Mathias Humbert, Pascal Berrang, Mario Fritz, and Michael Backes. MI-leaks: Model and data independent membership inference attacks and defenses on machine learning models. In *NCSS*, 2019.
- Sriram Sankararaman, Guillaume Obozinski, Michael I. Jordan, and Eran Halperin. Genomic privacy and limits of individual detection in a pool. *Nature Genetics*, 2009.
- Reza Shokri, Marco Stronati, and Vitaly Shmatikov. Membership inference attacks against machine learning models. *IEEE Symp. Security and Privacy*, 2017.
- Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2014.
- A. W. van der Vaart. *Asymptotic statistics*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 1998.
- Yu-Xiang Wang, Stephen Fienberg, and Alex Smola. Privacy for free: Posterior sampling and stochastic gradient monte carlo. In *ICML*, 2015.
- Yu-Xiang Wang, Jing Lei, and Stephen E Fienberg. On-average KL-privacy and its equivalence to generalization for max-entropy mechanisms. In *PSD*. Springer, 2016.

Max Welling and Yee Whye Teh. Bayesian learning via stochastic gradient langevin dynamics. In *ICML*, 2011.

Samuel Yeom, Irene Giacomelli, Matt Fredrikson, and Somesh Jha. Privacy risk in machine learning: Analyzing the connection to overfitting. In *CSF*, 2018.

Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. In *ICLR*, 2017.