# FUNDAMENTAL LIMITS OF MEMBERSHIP INFERENCE ATTACKS ON MACHINE LEARNING MODELS

**Eric Aubinais & Elisabeth Gassiat**
Université Paris-Saclay, CNRS
Laboratoire de mathématiques d'Orsay
Orsay, France
{eric.aubinais,elisabeth.gassiat}@universite-paris-saclay.fr

**Pablo Piantanida**
ILLS - International Laboratory on Learning Systems
MILA - Quebec AI Institute
CNRS, CentraleSupélec
Montreal, QC, Canada
pablo.piantanida@cnrs.fr

June 12, 2024

## ABSTRACT

Membership inference attacks (MIA) can reveal whether a particular data point was part of the training dataset, potentially exposing sensitive information about individuals. This article provides theoretical guarantees by exploring the fundamental statistical limitations associated with MIAs on machine learning models. More precisely, we first derive the statistical quantity that governs the effectiveness and success of such attacks. We then theoretically prove that in a non-linear regression setting with overfitting algorithms, attacks may have a high probability of success. Finally, we investigate several situations for which we provide bounds on this quantity of interest. Interestingly, our findings indicate that discretizing the data might enhance the algorithm's security. Specifically, it is demonstrated to be limited by a constant, which quantifies the diversity of the underlying data distribution. We illustrate those results through two simple simulations.

## 1 Introduction

In today's data-driven era, machine learning models are designed to reach higher performance, and the size of new models will then inherently increase, therefore the information stored (or memorized) in the parameters [Hartley and Tsaftaris, 2022, Del Grosso et al., 2023]. The protection of sensitive information is of paramount importance. Membership Inference Attacks (MIAs) have emerged as a concerning threat, capable of unveiling whether a specific data point was part of the training dataset of a machine learning model [Shokri et al., 2017, Song et al., 2017a, Nasr et al., 2019, Zhu et al., 2019]. Such attacks can potentially compromise individual privacy and security by exposing sensitive information [Carlini et al., 2023a]. Furthermore, the recent publication by Tabassi et al. [2019] from the National Institute of Standards and Technology (NIST) explicitly notes that an MIA that successfully identifies an individual as part of the dataset used for training the target model constitutes a breach of confidentiality.

To date, the most comprehensive defense mechanism against privacy attacks is Differential Privacy (DP), a framework initially introduced by Dwork et al. [2006]. DP has shown remarkable adaptability in safeguarding the privacy of machine learning models during training, as demonstrated by the works of Jayaraman and Evans [2019], Hannun et al. [2021]. However, it is worth noting that achieving a high level of privacy through differentially private training often comes at a significant cost to the accuracy of the model, especially when aiming for a low privacy parameter [Sablayrolles et al., 2019]. Conversely, when evaluating the practical effectiveness of DP in terms of its ability to protect against

privacy attacks empirically, the outlook is considerably more positive. DP has demonstrated its efficacy across a diverse spectrum of attacks, encompassing MIAs, attribute inference, and data reconstruction (see Guo et al. [2023] and references therein). DP has been extensively used to understand the performances of MIAs against learning systems Thudi et al. [2022] or how a mechanism could be introduced to defend oneself against MIAs He et al. [2022], Izzo et al. [2022].

Empirical evidence suggests that small models compared to the size of training set are often sufficient to thwart the majority of existent threats and empirically summarized in Baluta et al. [2022]. Similarly, when the architecture of a machine learning model is overcomplex with respect to the size of the training set, model overfitting increases the effectiveness of MIAs, as has been identified by Shokri et al. [2017], Yeom et al. [2018], He et al. [2022]. However, despite these empirical findings, there remains a significant gap in our theoretical understanding of this phenomenon. This article delves into the core statistical limitations surrounding MIAs on machine learning models at large.

Our investigation commences by establishing the **fundamental statistical quantity that governs the effectiveness and success of MIA attacks.** In the learning model under consideration, our focus lies on algorithms that can be described as a function of the empirical distribution of their training dataset.

Specifically, we concentrate on datasets of independent and identically distributed (*i.i.d.*) samples. To assess the effectiveness of MIAs, we will gauge their **accuracy** by examining their success probability in determining membership. Notably, we assess the security of a model based on the highest level of accuracy achieved among MIAs.

We delve into the intricacies of MIA and derive insights into the key factors that influence its outcomes. Subsequently, we explore various scenarios: overfitting algorithms, empirical mean-based algorithms and discrete data, among others, presenting bounds on this pivotal statistical quantity.

## 1.1 Contributions

In our research, we make theoretical contributions to the understanding of MIAs on machine learning models. Our key contributions can be summarized as follows:

- **Identification of Crucial Statistical Quantity:** We introduce the critical statistical quantity denoted as $\Delta_n(P, \mathcal{A})$, where $n$ represents the size of the training dataset, $P$ is the data distribution, and $\mathcal{A}$ is the underlying algorithm. This quantity plays a pivotal role in assessing the accuracy of effective MIAs. The quantity $\Delta_n(P, \mathcal{A})$ provides an intuitive measure of how distinct parameters of a model can be with respect to a sample in the training set, and as a result, it indicates the extent to which we can potentially recover sample membership through MIAs. Consequently, we demonstrate that when $\Delta_n(P, \mathcal{A})$ is small, the accuracy of the best MIA is notably constrained. Conversely, when $\Delta_n(P, \mathcal{A})$ approaches 1, the best MIA is successful with high probability. This highlights the importance of $\Delta_n(P, \mathcal{A})$ in characterizing information disclosure in relation to the training set.

- **Lower Bounds for Overfitting Algorithms**: For algorithms that overfit with high probability, we exhibit a lower bound on $\Delta_n(P, \mathcal{A})$ (see Theorem 4.3). In a non-linear regression setting, we further theoretically demonstrate that algorithms for which small training loss is reached, loss-based MIAs can achieve almost perfect inference, as illustrated in Section 6 by numerical experiments. Up to our knowledge, this is the first theoretical proof that overfitting indeed opens the way to successful MIAs.

- **Precise Upper Bounds for Empirical Mean-Based Algorithms:** For algorithms that compute functions of empirical means, we establish precise upper bounds on $\Delta_n(P, \mathcal{A})$. We prove that $\Delta_n(P, \mathcal{A})$ is bounded from above by a constant, determined by $(P, \mathcal{A})$, multiplied by $n^{-1/2}$. In practical terms, this means that having $\Omega(\varepsilon^{-2})$ samples in the dataset is sufficient to ensure that $\Delta_n(P, \mathcal{A})$ remains below $\varepsilon$ for any $\varepsilon \in (0, 1)$.

- **Maximization of $\Delta_n(P, \mathcal{A})$:** In scenarios involving discrete data with an infinite parameter space, we provide a precise formula for maximizing $\Delta_n(P, \mathcal{A})$ across all algorithms $\mathcal{A}$. Additionally, under specific assumptions, we determine that this maximization is proportional to $n^{-1/2}$ and to a quantity $C(P)$ which measures the diversity of the underlying data distribution. We illustrate this behaviour with numerical experiments in Section 6.

## 1.2 Related Works

**Privacy Attacks.** The majority of cutting-edge attacks follow a consistent approach within a framework known as Black-Box. In this framework, where access to the data distribution is available, attacks assess the performance of a model by comparing it to a group of "shadow models". These shadow models are trained with the same architecture but on an artificially and independently generated dataset from the same data distribution. Notably, loss evaluated on

training samples are expected to be much lower than when evaluated on "test points". Therefore, a significant disparity between these losses indicates that the sample in question was encountered during the training, effectively identifying it as a member. This is intuitively related to some sort of "stability" of the algorithm on training samples [Bousquet and Elisseeff, 2002]. Interestingly, we explicitly identify the exact quantity controlling the accuracy of effective MIAs which may be interpreted as a measure of stability of the underlying algorithm. In fact, as highlighted by Rezaei and Liu [2021], it is important to note that MIAs are not universally effective and their success depends on various factors. These factors include the characteristics of the data distribution, the architecture of the model, particularly its size, the size of the training dataset, and others, as discussed recently by Shokri et al. [2017], Carlini et al. [2022a]. Subsequently, there has been a growing body of research delving into Membership Inference Attacks (MIAs) on a wide array of machine learning models, encompassing regression models [Gupta et al., 2021], generation models [Hayes et al., 2018], and embedding models [Song and Raghunathan, 2020]. A comprehensive overview of the existing body of work on various MIAs has been systematically compiled in a thorough survey conducted by Hu et al. [2022]. While studies of MIAs through DP already reveal precise bounds, it is worth noting that these induce a significant loss of performance on the learning task. Interestingly, the findings of the Section 5 reveal a threshold on the minimum number of training samples to overcome the need of introducing DP mechanisms.

**Overfitting Effects.** The pioneering work by Shokri et al. [2017] has effectively elucidated the relationship between overfitting and the privacy risks inherent in many widely-used machine learning algorithms. These empirical studies clearly point out that overfitting can often provide attackers with the means to carry out membership inference attacks. This connection is extensively elaborated upon by Salem et al. [2018], Yeom et al. [2018], and later by He et al. [2022], among other researchers. Overfitting tends to occur when the underlying model has a complex architecture or when there is limited training data available, as explained in Baluta et al. [2022]. Recent works [Yeom et al., 2018, Del Grosso et al., 2023] investigated the theoretical aspects of the overfitting effect on the performances of MIAs, showing that the MIA performances can be lower bounded by a function of the *generalization gap* under some assumptions on the loss function. In our paper, we explicitly emphasize these insights by quantifying the dependence of $\Delta_n(P, \mathcal{A})$ either on the dataset size and underlying structural parameters, or explicitly on the overfitting probability of the learning model.

**Memorization Effects.** Machine learning models trained on private datasets may inadvertently reveal sensitive data due to the nature of the training process. This potential disclosure of sensitive information occurs as a result of various factors inherent to the training procedure, which include the extraction of patterns, associations, and subtle correlations from the data [Song et al., 2017a, Zhang et al., 2021]. While the primary objective is to generalize from data and make predictions, there is a risk that these models may also pick up on, and inadvertently expose, confidential or private information contained within the training data. This phenomenon is particularly concerning as it can lead to privacy breaches, compromising the confidentiality and security of personal or sensitive data [Hartley and Tsaftaris, 2022, Carlini et al., 2022b, 2019, Leino and Fredrikson, 2020, Thomas et al., 2020]. Recent empirical studies have shed light on the fact that, in these scenarios, it is relatively rare for the average data point to be revealed by learning models [Tirumala et al., 2022, Murakonda and Shokri, 2007, Song et al., 2017b]. What these studies have consistently shown is that it is the outlier samples that are more likely to undergo memorization by the model [Feldman, 2020], leading to potential data leakage. This pattern can be attributed to the nature of learning algorithms, which strive to generalize from the data and make predictions based on common patterns and trends. Average or typical data points tend to conform to these patterns and are thus less likely to stand out. On the other hand, outlier samples, by their very definition, deviate significantly from the norm and may capture the attention of the model. So when an outlier sample is memorized, it means the model has learned it exceptionally well, potentially retaining the unique characteristics of that data point. As a consequence, when exposed to similar data points during inference, the model may inadvertently leak information it learned from the outliers, compromising the privacy and security of the underlying data. An increasing body of research is dedicated to the understanding of memorization effects in language models [Carlini et al., 2023b]. In the context of our research, it is important to highlight that our primary focus is on understanding the accuracy of MIAs but not its relationship with memorization. Indeed, this connection remains an area of ongoing exploration and inquiry in our work.

## 2 Background and Problem Setup

In this paper, we focus on MIAs, the ability of recovering membership to a training dataset $\mathbf{z} := (z_1, \cdots, z_n) \in \mathcal{Z}^n$ of a test point $\tilde{z} \in \mathcal{Z}$ from a predictor $\hat{\mu} = \mu_{\hat{\theta}_n}$ in a model $\mathcal{F} := \{\mu_\theta : \theta \in \Theta\}$, where $\Theta$ is the space of parameters. The predictor is identified to its parameters $\hat{\theta}_n \in \Theta$ learned from $\mathbf{z}$ through an **algorithm** $\mathcal{A} : \bigcup_{k>0} \mathcal{Z}^k \to \mathcal{P}' \subseteq \mathcal{P}(\Theta)$, that is $\hat{\theta}_n$ follows the distribution $\mathcal{A}(\mathbf{z})$ conditionally to $\mathbf{z}$, which we assume we have access to. Here, $\mathcal{P}(\Theta)$ is the set of all distributions on $\Theta$, and $\mathcal{P}'$ is the range of $\mathcal{A}$.

This means that there exists a function $g$ and a random variable $\xi$ independent of $\mathbf{z}$ such that $\hat{\theta}_n = g(\mathbf{z}, \xi)$. When $\mathcal{A}$ takes values in the set of Dirac distributions, that is $\hat{\theta}_n$ is a deterministic function of the data, we shall identify the parameters directly to the output of the algorithm $\hat{\theta}_n \coloneqq \mathcal{A}(\mathbf{z}_1, \cdots, \mathbf{z}_n)$.

Throughout the paper, we will further assume that the algorithm $\mathcal{A}$ can be expressed as a function of the empirical distribution of the training dataset. Letting $\mathcal{M}$ be the set of all discrete distributions on $\mathcal{Z}$, and $\hat{P}_n$ be the empirical distribution of the training dataset, it means that there exists a (randomized) function $G : \mathcal{M} \to \mathcal{P}'$ such that we have $\mathcal{A}(\mathbf{z}_1, \cdots, \mathbf{z}_n) = G(\hat{P}_n)$ (almost surely).

Interestingly, if an algorithm minimizes an empirical cost, then it satisfies this assumption. In particular, maximum likelihood based algorithms or Bayesian methods from Sablayrolles et al. [2019] are special cases. Any instance of an algorithm in what follows will satisfy these assumptions.

We further discuss this assumption in Appendix A.

We consider MIAs as functions of the parameters and the test point whose outputs are 0 or 1.

**Definition 2.1** (Membership Inference Attack - MIA). *Any measurable map $\phi : \Theta \times \mathcal{Z} \to \{0, 1\}$ is called a **Membership Inference Attack**.*

In varying contexts, MIAs might access more information, including randomization. While we omit these extra details here, it's worth noting that our results remain applicable.

We measure the accuracy of an MIA $\phi$ by its probability of successfully guessing the membership of the test point. For that purpose, we encode membership to the training data set as 1. We assume that $\mathbf{z}_1, \ldots, \mathbf{z}_n$ are independent and identically distributed (*i.i.d.*) random variables with distribution $P$. Following Del Grosso et al. [2023] or Sablayrolles et al. [2019] framework, we suppose that the test point $\tilde{z}$ is to be drawn from $P$ independently from the samples $\mathbf{z}_1, \ldots, \mathbf{z}_n$ with probability $\nu \in (0, 1)$. Otherwise, conditionally to $\mathbf{z}$, we set $\tilde{z}$ to any $\mathbf{z}_j$ each with uniform probability $1/n$.

Letting $U$ be a random variable with distribution $\hat{P}_n \coloneqq \frac{1}{n} \sum_{j=1}^{n} \delta_{\mathbf{z}_j}$ conditionally to $\mathbf{z}$, $\mathbf{z}_0$ to be drawn independently from $P$ and $T$ be a random variable having Bernoulli distribution with parameter $\nu$ and independent of any other random variables, we can state

$$\tilde{z} \coloneqq T\mathbf{z}_0 + (1 - T)U.$$

**Definition 2.2** (Accuracy of an MIA). *The **accuracy of an MIA** $\phi$ is defined as*

$$Acc_n(\phi; P, \mathcal{A}) \coloneqq P\left(\phi(\hat{\theta}_n, \tilde{z}) = 1 - T\right), \tag{1}$$

*where the probability is taken over all randomness.*

The accuracy of an MIA scales from 0 to 1. Constant MIAs $\phi_0 \equiv 0$ and $\phi_1 \equiv 1$ have respectively an accuracy equal to $\nu$ and $1 - \nu$, which means that we always can build an MIA with accuracy of at least $\max(\nu, 1 - \nu)$ and any MIA performing worse than this quantity is irrelevant to use. We now define the **Membership Inference Security** of an algorithm as a quantity summarizing the amount of security of the system against MIAs[1].

**Definition 2.3** (Membership Inference Security - MIS). *Let $\nu_* \coloneqq \min(\nu, 1 - \nu)$. The **Membership Inference Security** of an algorithm $\mathcal{A}$ is*

$$Sec_n(P, \mathcal{A}) \coloneqq \nu_*^{-1}\left(1 - \sup_\phi Acc_n(\phi; P, \mathcal{A})\right), \tag{2}$$

*where the sup is taken over all MIAs.*

The Membership Inference Security scales from 0 (the best MIA approaches perfect guess of membership) to 1 (MIAs can not do better than $\phi_0$ and $\phi_1$).

## 3   Performance Assessment of Membership Inference Attacks

In this section, we prove that the **Crucial Statistical Quantity** for the assessment of the accuracy of membership inference attacks is $\Delta_n(P, \mathcal{A})$, defined as

$$\Delta_n(P, \mathcal{A}) \coloneqq \left\|\mathcal{L}\left((\hat{\theta}_n, \mathbf{z}_1)\right) - \mathcal{L}\left((\hat{\theta}_n, \mathbf{z}_0)\right)\right\|_{\mathrm{TV}}, \tag{3}$$

---

[1]Depending on the application, the True Positive Rate (TPR) and the False Positive Rate (FPR) provide a more flexible and customizable approach to evaluating attacker performance, especially in scenarios where the cost of false positives and false negatives may differ. The extension of our results to this framework is relegated to future work.

which depends on $P$, $n$ and $\mathcal{A}$. Here, for any random variable x, $\mathcal{L}(x)$ denotes its probability distribution, and for any distributions $Q_1$ and $Q_2$, $\|Q_1 - Q_2\|_{\mathrm{TV}}$ denotes the total variation distance between $Q_1$ and $Q_2$. One can interpret $\Delta_n(P, \mathcal{A})$ as quantifying some stability of the algorithm. As per the $i.i.d.$ assumption of the data, the choice of $z_1$ is arbitrary and is only for simplicity purpose.

**Theorem 3.1** (Key bound on accuracy). *Suppose $P$ is any distribution and $\mathcal{A}$ is any algorithm. Then the accuracy of any MIA $\phi$ satisfies:*

$$\nu_* - \nu_* \Delta_n(P, \mathcal{A}) \leq Acc_n(\phi; P, \mathcal{A}) \leq 1 - \nu_* + \nu_* \Delta_n(P, \mathcal{A}).$$

*In particular,*

$$Sec_n(P, \mathcal{A}) \geq 1 - \Delta_n(P, \mathcal{A}),$$

*with equality when $\nu = 1/2$.*

We see that $\Delta_n(P, \mathcal{A})$ appears to be a key mathematical quantity for assessing the accuracy of MIAs. Theorem 3.1 shows that an upper bound on $\Delta_n(P, \mathcal{A})$ translates into a lower bound for the MIS of any algorithm; and when $\nu = 1/2$, $\Delta_n(P, \mathcal{A})$ is the quantity that controls the best possible accuracy of MIAs. We thus study in Section 4 a control on $\Delta_n(P, \mathcal{A})$ when the algorithm overfits, and in Section 5 situations in which we are able to give precise controls on $\Delta_n(P, \mathcal{A})$. We give in Section 6 some numerical experiments. Proof of the Theorem can be found in Appendix E.

**Remark :** Notice that there is no assumption on the data distribution $P$. For instance, we can take into account outliers by making $P$ a mixture.

## 4 Overfitting Causes Lack of Security

In this section, we assume that $\mathcal{Z} \coloneqq \mathcal{X} \times \mathcal{Y}$ and that the algorithm $\mathcal{A}$ produces overfitting parameters $\hat{\theta}_n$. We then note $z \coloneqq (x, y)$. We consider learning systems minimizing $L_n : \theta \mapsto \frac{1}{n} \sum_{i=1}^{n} l_\theta(x_i, y_i)$ for some training dataset $(z_1, \cdots, z_n)$ where $l_\theta : \mathcal{X} \times \mathcal{Y} \to \mathbb{R}^+$ is a loss function. We defer all proofs of this section to Appendix F.

**Definition 4.1** (($\varepsilon, 1 - \alpha$)-Overfitting). *We say that the algorithm $\mathcal{A}$ is ($\varepsilon, 1 - \alpha$)-overfitting for some $\varepsilon \in \mathbb{R}^+$ and $\alpha \in (0, 1)$ when*

$$P\left(l_{\hat{\theta}_n}(x_1, y_1) \leq \varepsilon\right) \geq 1 - \alpha. \tag{4}$$

When $\alpha = 0$, Equation 4 is equivalent to having $l_{\hat{\theta}_n}(x_i, y_i) \leq \varepsilon$ almost surely for all $i = 1, \ldots, n$. Furthermore, in many algorithms, we give an additional stopping criteria taking the form $L_n \leq \eta$ for some $\eta \in \mathbb{R}^+$. Letting $\mathcal{A}_\eta$ such an algorithm, we give a sufficient condition for Equation 4 to hold:

**Proposition 4.2.** *For some fixed $\varepsilon \in \mathbb{R}^+$ and $\alpha \in (0, 1)$, let $\eta \coloneqq \varepsilon\alpha$ and suppose that $\mathcal{A}_\eta$ stops as soon as $L_n(\hat{\theta}_n) \leq \eta$. Then $\mathcal{A}_\eta$ is ($\varepsilon, 1 - \alpha$)-overfitting.*

We will need the an additional hypothesis for the following theorem.

**Hypothesis (H1) :** $\mathcal{Y} \coloneqq \mathbb{R}^s$ for some $s \geq 1$ for all $\theta \in \Theta$, $x \in \mathcal{X}$ and $y \in \mathcal{Y}$, we have

$$l_\theta(x, y) = \omega(y, \Psi_\theta(x)), \tag{5}$$

for some family of functions $\Psi_\theta : \mathcal{X} \to \mathbb{R}$ and some continuous function $\omega : \mathbb{R} \times \mathbb{R} \to \mathbb{R}$.

**Theorem 4.3** (Overfitting induces lack of security). *Assume $\mathcal{A}$ is ($\varepsilon, 1 - \alpha$)-overfitting for some fixed ($\varepsilon, \alpha$). Let $S_\theta^\varepsilon \coloneqq \{(x, y) \in \mathcal{X} \times \mathcal{Y} : l_\theta(x, y) \leq \varepsilon\}$ be the $\varepsilon$-sub-level set of $l_\theta$ for all $\theta \in \Theta$. Then we have*

$$\Delta_n(P, \mathcal{A}) \geq 1 - \alpha - \int_\Theta P((x, y) \in S_\theta^\varepsilon) d\mu_{\hat{\theta}_n}, \tag{6}$$

*where $\mu_{\hat{\theta}_n}$ is the distribution of $\hat{\theta}_n$.*

*Assume furthermore that **H1** holds. Assume that for all $\eta > 0$, $\mathcal{A}_\eta$ stops as soon as $L_n \leq \eta$ and that a version of the conditional distribution of y given x is absolutely continuous with respect to the Lebesgue measure, then*

$$\lim_{\eta \to 0^+} \Delta_n(P, \mathcal{A}_\eta) = 1. \tag{7}$$

The second point of Theorem 4.3 states that for regressors with reasonably low training loss on the dataset, a loss-based MIA would reach high success probability. This theoretically confirms the already well-known insight that overfitting implies poor security. We further discuss these hypotheses in Appendix B.

Hypothesis **H1** occurs when $\Psi_\theta(x)$ models the conditional expectation of y given x, in a setting where the loss function is defined as a distance between $\Psi_\theta(x)$ and y.

**Example 4.4** (non-linear regression Neural Network). *We consider here a (non-linear) regression setting, that is for all $j = 1, \ldots, n$, we have $y_j := \Psi^*(\mathrm{x}_j) + \zeta_j$, where $\zeta_j$ is some independent random noise and the function $\Psi^* : \mathcal{X} \to \mathbb{R}$ is arbitrary, fixed and unknown. We aim at estimating $\Psi^*$ by some Neural Network $\Psi_\theta \in \mathcal{F}$, where $\mathcal{F}$ is some fixed model. For instance $\mathcal{F}$ can be the set of all 2-layers ReLU neural networks with fixed hidden layer width. The algorithm $\mathcal{A}$ then learns by minimizing the MSE loss $L_n := \frac{1}{n} \sum_{j=1}^n (y_j - \Psi_\theta(x_j))^2$. In this case, Equation 5 holds. Under the further assumption that there is an arbitrarily close approximation $\Psi_\theta$ of $\Psi^*$ in $\mathcal{F}$, one can construct the sequence of algorithms $(\mathcal{A}_\eta)_{\eta \in \mathbb{R}^+}$ such that the hypotheses of the second point of Theorem 4.3 for Equation 7 to hold. Refer to Section 6 for a numerical illustration.*

**Example 4.5** (Linear regression). *We assume here a linear regression setting, that is $\mathcal{X} := \mathbb{R}^d$ for some $d \in \mathbb{N}$, and $y_j := \beta^T \mathrm{x}_j + \zeta_j$, where $\zeta_j$ is some independent random noise and $\beta \in \mathbb{R}^d$ is fixed and unknown. Further assuming that $\zeta_j$ is absolutely continuous with respect to the Lebesgue measure, and that $d > n$, both Equations 4 (with $\varepsilon, \alpha = 0$) and 5 hold. Then, the assumptions of the second point of Theorem 4.3 are satisfied, leading to $\Delta_n(P, \mathcal{A}) = 1$.*

# 5 Security is Data Size Dependent

In this section, we study the converse, where we aim at understanding when to expect $\Delta_n(P, \mathcal{A})$ to be close to 0. All the proofs of the section can be found in Appendix G.

## 5.1 Empirical Mean based algorithms

We first study the case of algorithms for which the parameters $\hat{\theta}_n$ can be expressed in the form of functions of empirical means (e.g., linear regression with mean-squared error, method of moments...). Specifically, for any (fixed) measurable maps $L : \mathcal{Z} \to \mathbb{R}^d$ and $F : \mathbb{R}^d \to \mathbb{R}^q$ for some $d, q \in \mathbb{N}$, we consider that

$$\hat{\theta}_n := F\left(\frac{1}{n} \sum_{j=1}^n L(z_j)\right). \tag{8}$$

Equation 8 states that the parameters are the result of the algorithm $\mathcal{A} : (z_1, \cdots, z_n) \mapsto \delta_{F\left(\frac{1}{n} \sum_{j=1}^n L(z_j)\right)}$, where $\delta_\theta$ stands for the Dirac mass at $\theta$. We then have the following result

**Theorem 5.1.** *Suppose that the distribution of $L(\mathrm{z}_1)$ has a non zero absolutely continuous part with respect to the Lebesgue measure, and a third finite moment. Then*

$$\Delta_n(P, \mathcal{A}) \le c_{L,P} n^{-1/2} + \frac{\sqrt{d}}{2n}, \tag{9}$$

*for some constant $c_{L,P}$ depending only on $L$ and $P$.*

**Remark 5.1** : Theorem 5.1 implies that a sufficient condition to ensure $Sec_n(P, \mathcal{A})$ to be made larger than $1 - \varepsilon$, is to have $\Delta_n(P, \mathcal{A}) \le \varepsilon$ which holds as soons as $n \ge \Omega(\varepsilon^{-2})$. The hidden constant only depends on the distribution data $P$ and the parameters dimension $d$. See Appendix G for a proof.

We now provide examples for which Theorem 5.1 allows us to give an upper bound on $\Delta_n(P, \mathcal{A})$.

**Example 5.2** (solving equations). *We seek to estimate an (unknown) parameter of interest $\theta_0 \in \Theta \subseteq \mathbb{R}^d$. We suppose that we are given two functions $h : \Theta \to \mathbb{R}^l$ and $\psi : \mathcal{Z} \to \mathbb{R}^l$ for some $l \in \mathbb{N}$, and that $\theta_0$ is solution to the equation*

$$h(\theta_0) = \mathbb{E}[\psi(\mathrm{z})]. \tag{10}$$

*where $\mathrm{z}$ is a random variable of distribution $P$. Having access to data samples $\mathrm{z}_1, \ldots, \mathrm{z}_n$ drawn independently from the distribution $P$, we estimate $\mathbb{E}[\psi(\mathrm{z})]$ by $\frac{1}{n} \sum_{j=1}^n \psi(\mathrm{z}_j)$. Assuming that $h$ is invertible, one can set $\hat{\theta}_n = h^{-1}\left(\frac{1}{n} \sum_{j=1}^n \psi(\mathrm{z}_j)\right)$, provided that $\frac{1}{n} \sum_{j=1}^n \psi(\mathrm{z}_j) \in \mathcal{I}m(h)$. In particular, when $\mathcal{Z} = \mathbb{R}$, this method generalizes the method of moments by setting $\psi(z) = (z, z^2, \cdots, z^l)$. We then may apply Theorem 5.1 to any estimators obtained by solving equations.*

**Example 5.3** (Linear Regression). *We consider here the same framework as in Example 4.5, where $d < n$ (hence Definition 4.1 can not be fulfilled with $\alpha = 0$). Let $\mathbb{X}$ be the $d \times n$ matrix whose $i^{th}$ row is $\mathrm{x}_i$, and $\mathbb{Y}$ be the column vector $(\mathrm{y}_1, \cdots, \mathrm{y}_n)^T$. We then recall that the estimator $\hat{\beta}_n$ of $\beta$ is given by*

$$\hat{\beta}_n := (\mathbb{X}\mathbb{X}^T)^{-1} \mathbb{X}\mathbb{Y}^T.$$

*Based on Equation 8, if we set $F(K, b) := K^{-1}b^T$ and $L((x, y)) := ((x^i x^j)_{i,j=1}^d, (x^i y)_{i=1}^d)$, where $x^i$ is the $i^{th}$ coordinate of $x$, then we can express the estimator as $\hat{\beta}_n = F\left(\frac{1}{n}\sum_{j=1}^n L((x_j, y_j))\right)$.*

Interestingly, we see from Examples 4.5 and 5.3 that the security of least squares linear regression estimator is constant $0$ up to $n = d$ (where $d$ is both the dimension of the data and the dimension of the parameters), and then is increasing up to $1$ provided that $n \to \infty$.

## 5.2 Discrete Data Distribution

We now consider the common distribution of the points in the data set to be $P := \sum_{j=1}^K p_j \delta_{u_j}$ for some fixed $K \in \mathbb{N} \cup \{\infty\}$, some fixed probability vector $(p_1, \cdots, p_K)$ and some fixed points $u_1, \ldots, u_K$ in $\mathcal{Z}$. Without loss of generality, we may assume that $p_j > 0$ for all $j \in \{1, \cdots, K\}$.

**Theorem 5.4.** *For $j = 1, \ldots, K$, let $B_j$ be a random variable having Binomial distribution with parameters $(n, p_j)$. Then,*

$$\max_{\mathcal{A}} \Delta_n(P, \mathcal{A}) = \frac{1}{2}\sum_{j=1}^K \mathbb{E}\left[\left|\frac{B_j}{n} - p_j\right|\right], \tag{11}$$

*where the $\max$ is taken over all algorithms and is reached on algorithms of the form $\mathcal{A}(z_1, \cdots, z_n) = \delta_{F(\frac{1}{n}\sum_{j=1}^n \delta_{z_j})}$ for some injective maps $F$.*

Theorem 5.4 gives a precise formula to accurately bound $\Delta_n(P, \mathcal{A})$ for any algorithm $\mathcal{A}$. We show below that the r.h.s. of Equation 11 is tightly related to the quantity $C(P) := \sum_{j=1}^K \sqrt{p_j(1-p_j)}$. It is worth noting that $C(P)$ is a diversity measure, giving a control on the homogeneity of the data distribution. We show in Appendix C that it is comparable both to the Gini-Simpson and the Shannon Entropy.

**Corollary 5.4.1.** *Assume that $C(P) < \infty$, $n \geq 5$ and $n > 1/p_j$ for all $j = 1, \ldots, n$. Then there exists a universal constant $c \geq 0.29$ such that*

$$c \cdot C(P)n^{-1/2} \leq \max_{\mathcal{A}} \Delta_n(P, \mathcal{A}) \leq \frac{C(P)}{2}n^{-1/2}.$$

*If $C(P) < \infty$ but the condition on $n$ does not hold, we still have $\max_{\mathcal{A}} \Delta_n(P, \mathcal{A}) \leq \frac{C(P)}{2}n^{-1/2}$.*

Corollary 5.4.1 implies that a sufficient condition to ensure security larger than $1 - \varepsilon$ is to have at $n \geq (C(P)/2\varepsilon)^2$. In any case, this result indicates that discretizing the data allows the control on the security of any algorithm. Section 6 illustrates the impact of $C(P)$ through numerical experiments. We further discuss this Corollary in Section C.

# 6 Numerical Experiments

In this section, we propose two numerical experiments to illustrate our results in Sections 4 and 5. All simulations have been conducted with Pytorch library. We refer to Appendix D for more details on the experiments.

## 6.1 Overfitting

We run a non-linear regression experiment to illustrate the results of Section 4 and specifically Example 4.4. We then consider the setting of Example 4.4 with $\Psi^*(x) = \sin(\pi\beta^T x)$ for some fixed $\beta$. During the training of the neural network $\Psi_{\hat{\theta}_n}$, at each iteration we evaluate the fraction of training (validation) data that achieves a loss below $\varepsilon$. Validation data correspond to a set of data independent from the training dataset.

Figure 1 illustrates Theorem 4.3 by showing that for very small values of threshold ($\varepsilon = 10^{-6}$), we still reach $100\%$ training accuracy after 2500 iterations whereas the validation accuracy (for $\varepsilon = 10^{-6}$) stabilizes at near $0\%$. In this case, a simple loss-based MIA with threshold $\varepsilon$ would suffice to accurately predict membership most of the case. The number of iterations being generally unknown to the MIA, the calibration of $\varepsilon$ is a hard task to perform. In Figure 1, even though it seems that the loss-based attack with threshold $\varepsilon = 10^{-6}$ is a good candidate to achieve near perfect guess, it is worth noting that it would occur only if at least 2000 iterations have been done during the training procedure.
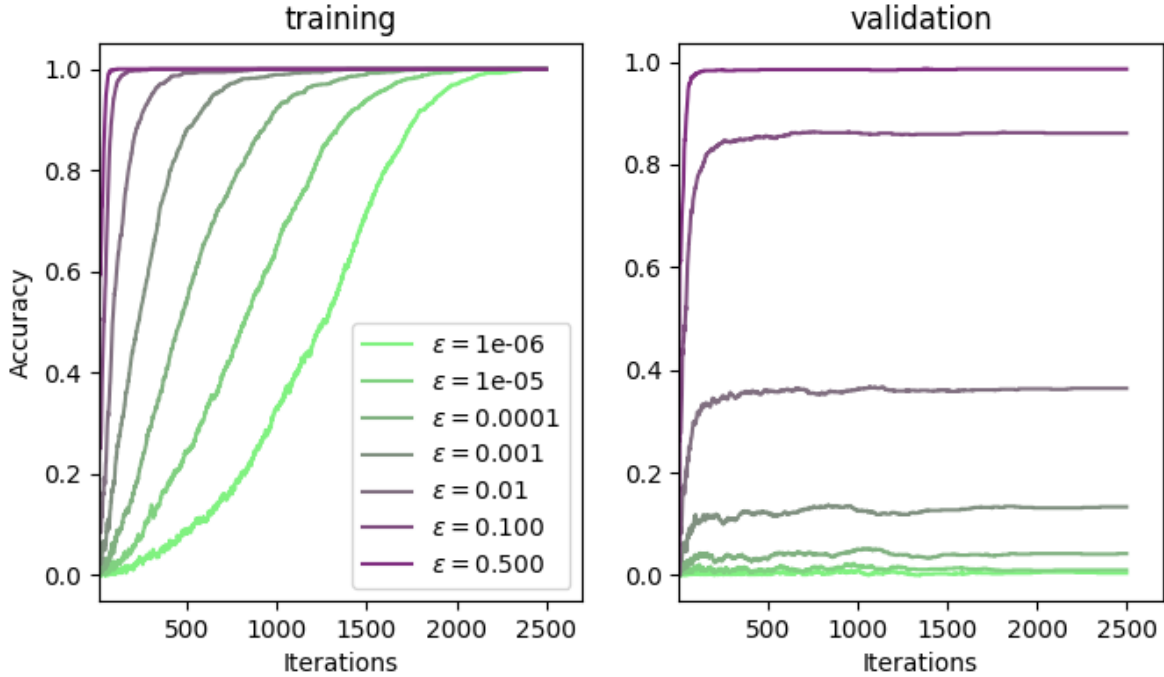
Figure 1: Shows the fraction of the Training/Validation dataset whose loss is under given thresholds during the training process. The left figure shows the **training accuracy**, and the right figure shows the **validation accuracy**.

## 6.2 Impact of $C(P)$ on accuracy

Corollary 5.4.1 indicates that discretizing the data distribution improves the security of the model. To illustrate the impact of a discretization through the constant $C(P)$, we trained several 3-layers neural networks to classify samples from the MNIST dataset [Deng, 2012]. Before training, we fixed three discretizations[2] (clusterings). For each dataset (with varying size $n$), we trained a neural network on it, and three other neural networks on discretized versions of the original dataset (one for each clustering). We then numerically computed the quantity $C(P)$ for each discretization. Table 1 shows the accuracy of all neural networks, and displays the impact of the discretization on the accuracy,

| n | raw dataset | $C(P) = 4.3$ | $C(P) = 6.74$ | $C(P) = 9.20$ |
|---|---|---|---|---|
| 1000 | $0.989 \pm 0.0011$ | $0.963 \pm 0.0223 \, (\Delta_n \leq 0.07)$ | $0.968 \pm 0.0137 \, (\Delta_n \leq 0.11)$ | $0.986 \pm 0.0039 \, (\Delta_n \leq 0.15)$ |
| 5000 | $0.993 \pm 0.0012$ | $0.967 \pm 0.0184 \, (\Delta_n \leq 0.03)$ | $0.971 \pm 0.0282 \, (\Delta_n \leq 0.05)$ | $0.984 \pm 0.0055 \, (\Delta_n \leq 0.07)$ |
| 10000 | $0.994 \pm 0.0006$ | $0.971 \pm 0.0141 \, (\Delta_n \leq 0.02)$ | $0.977 \pm 0.0082 \, (\Delta_n \leq 0.03)$ | $0.984 \pm 0.0055 \, (\Delta_n \leq 0.05)$ |

Table 1: Shows the accuracy of classifiers on MNIST dataset. The column **n** displays the dataset size. The column **raw dataset** displays the accuracy of the neural network on the original dataset. Each column of the column **discretized datasets** displays the accuracy of a neural network on the discretized dataset associated to the constant $C(P)$.

depending on $n$ and the value of $C(P)$. For a dataset of size $n = 1000$, our neural network reaches an accuracy of $0.989$ when trained on the original dataset. When discretizing, Table 1 displays a loss of almost $2.5\%$ of accuracy for the discretization having $C(P) = 4.30$, and a loss of $2\%$ for the other discretizations. As discussed in Section C, increasing the number of clusters will increase the value of $C(P)$. Table 1 displays the intuition that smaller discretization (smaller $C(P)$) will lower simultaneously the accuracy and the quantity $\Delta_n(P, \mathcal{A})$, which motivates the need to optimize the discretization to find a trade-off between security and accuracy.

---

[2]Many clustering algorithms exist, but we did not aim at optimizing the choice of the discretizations.

# 7 Summary and Discussion

The findings presented in this article open gates to the theoretical understanding of MIAs, and partially confirm some of empirically observed facts. Specifically, we confirmed that overfitting indeed induces the possibility of highly successful attacks. We further revealed a sufficient condition on the size of the training dataset to ensure control on the security of the learning algorithm, when dealing with discrete data distributions or functionals of empirical means. We established that the rates of convergence consistently follow an order of $n^{-1/2}$. The constants established in the rates of convergence scale with the number of discrete data points and the dimension of the parameters in the case of functionals of empirical means. Interestingly, for discrete data, the quantity $C(P)$ which is a diversity measure, highlights the use of data quantization to ensure privacy by design.

**Limitations and perspectives for further extensions of the present work.** In Section 5, our work is currently limited to the discrete data and empirical mean based algorithms. We intend to extend further our research to the complete study of maximum likelihood estimation, empirical loss minimization and Stochastic Gradient Descent. Additionally, we plan to expand our study to include quantized parameters. Furthermore, we aim to explore the optimization of the trade-off discussed at the conclusion of Section 6.2.

Our findings about overfitting algorithms do not cover classification algorithms. We anticipate continuing our research in this direction to gain a comprehensive understanding of the impact of overfitting. Currently, we are able to establish a result for $2-$ReLU binary classification networks, albeit under very stringent assumptions. Specifically, we study the case when the data are concentrated on the sphere of radius $\sqrt{s}$ in $\mathbb{R}^s$ and in a high-dimensional setting $s \geq n$. When the algorithm outputs a classifier whose parameters are in the direction of the gradient flow minimizing the exponential loss or the logistic loss, we prove that $\Delta_n(P, \mathcal{A})$ is lower bounded by the probability of the data to be not far from orthogonality, see Appendix B Proposition B.1. We anticipate that these assumptions may be relaxed in future investigations.

# References

John Hartley and Sotirios A Tsaftaris. Measuring unintended memorisation of unique private features in neural networks. *arXiv preprint arXiv:2202.08099*, 2022.

Ganesh Del Grosso, Georg Pichler, Catuscia Palamidessi, and Pablo Piantanida. Bounding information leakage in machine learning. *Neurocomputing*, 534:1–17, 2023.

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.

Congzheng Song, Thomas Ristenpart, and Vitaly Shmatikov. Machine Learning Models That Remember Too Much. In *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*, CCS '17, page 587–601, New York, NY, USA, 2017a. Association for Computing Machinery. ISBN 9781450349468. doi:10.1145/3133956.3134077. URL https://doi.org/10.1145/3133956.3134077.

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, 2019. doi:10.1109/SP.2019.00065.

Ligeng Zhu, Zhijian Liu, and Song Han. Deep leakage from gradients. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL https://proceedings.neurips.cc/paper_files/paper/2019/file/60a6c4002cc7b29142def8871531281a-Paper.pdf.

Nicolas Carlini, Jamie Hayes, Milad Nasr, Matthew Jagielski, Vikash Sehwag, Florian Tramer, Borja Balle, Daphne Ippolito, and Eric Wallace. Extracting training data from diffusion models. In *32nd USENIX Security Symposium (USENIX Security 23)*, pages 5253–5270, 2023a.

Elham Tabassi, Kevin Burns, Michael Hadjimichael, Andres Molina-Markham, and Julian Sexton. A taxonomy and terminology of adversarial machine learning, 10 2019.

Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity in private data analysis. In *Theory of Cryptography: Third Theory of Cryptography Conference, TCC 2006, New York, NY, USA, March 4-7, 2006. Proceedings 3*, pages 265–284. Springer, 2006.

Bargav Jayaraman and David Evans. Evaluating Differentially Private Machine Learning in Practice. In *Proceedings of the 28th USENIX Conference on Security Symposium*, SEC'19, page 1895–1912, USA, 2019. USENIX Association. ISBN 9781939133069.

Awni Y. Hannun, Chuan Guo, and Laurens van der Maaten. Measuring Data Leakage in Machine-Learning Models with Fisher Information. In *Conference on Uncertainty in Artificial Intelligence*, 2021. URL `https://api.semanticscholar.org/CorpusID:232013768`.

Alexandre Sablayrolles, Matthijs Douze, Cordelia Schmid, Yann Ollivier, and Hervé Jégou. White-box vs black-box: Bayes optimal strategies for membership inference. In *International Conference on Machine Learning*, pages 5558–5567. PMLR, 2019.

Chuan Guo, Alexandre Sablayrolles, and Maziar Sanjabi. Analyzing Privacy Leakage in Machine Learning via Multiple Hypothesis Testing: A Lesson From Fano. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett, editors, *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 11998–12011. PMLR, 23–29 Jul 2023. URL `https://proceedings.mlr.press/v202/guo23e.html`.

Anvith Thudi, Ilia Shumailov, Franziska Boenisch, and Nicolas Papernot. Bounding membership inference. *arXiv preprint arXiv:2202.12232*, 2022.

Xinlei He, Zheng Li, Weilin Xu, Cory Cornelius, and Yang Zhang. Membership-Doctor: Comprehensive Assessment of Membership Inference Against Machine Learning Models. *arXiv preprint arXiv:2208.10445*, 2022.

Zachary Izzo, Jinsung Yoon, Sercan O Arik, and James Zou. Provable Membership Inference Privacy. *arXiv preprint arXiv:2211.06582*, 2022.

Teodora Baluta, Shiqi Shen, S Hitarth, Shruti Tople, and Prateek Saxena. Membership inference attacks and generalization: A causal perspective. In *Proceedings of the 2022 ACM SIGSAC Conference on Computer and Communications Security*, pages 249–262, 2022.

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.

Olivier Bousquet and André Elisseeff. Stability and Generalization. *Journal of Machine Learning Research*, 2(Mar): 499–526, 2002. ISSN ISSN 1533-7928. URL `http://www.jmlr.org/papers/v2/bousquet02a.html`.

Shahbaz Rezaei and Xin Liu. On the difficulty of membership inference attacks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7892–7900, 2021.

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, 2022a.

Umang Gupta, Dmitris Stripelis, Pradeep Lam, Paul M. Thompson, J. Ambite, and Greg Ver Steeg. Membership Inference Attacks on Deep Regression Models for Neuroimaging. In *International Conference on Medical Imaging with Deep Learning*, 2021. URL `https://api.semanticscholar.org/CorpusID:233864706`.

Jamie Hayes, Luca Melis, George Danezis, and Emiliano De Cristofaro. Membership Inference Attacks Against Generative Models. 2018. URL `https://api.semanticscholar.org/CorpusID:202588705`.

Congzheng Song and Ananth Raghunathan. Information Leakage in Embedding Models. In *Proceedings of the 2020 ACM SIGSAC Conference on Computer and Communications Security*, CCS '20, page 377–390, New York, NY, USA, 2020. Association for Computing Machinery. ISBN 9781450370899. doi:10.1145/3372297.3417270. URL `https://doi.org/10.1145/3372297.3417270`.

Hongsheng Hu, Zoran Salcic, Lichao Sun, Gillian Dobbie, Philip S Yu, and Xuyun Zhang. Membership inference attacks on machine learning: A survey. *ACM Computing Surveys (CSUR)*, 54(11s):1–37, 2022.

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.

Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding Deep Learning (Still) Requires Rethinking Generalization. *Commun. ACM*, 64(3):107–115, feb 2021. ISSN 0001-0782. doi:10.1145/3446776. URL `https://doi.org/10.1145/3446776`.

Nicholas Carlini, Matthew Jagielski, Chiyuan Zhang, Nicolas Papernot, Andreas Terzis, and Florian Tramer. The Privacy Onion Effect: Memorization is Relative. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems*, volume 35, pages 13263–13276. Curran Associates, Inc., 2022b.

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.

Klas Leino and Matt Fredrikson. Stolen memories: Leveraging model memorization for calibrated {White-Box} membership inference. In *29th USENIX security symposium (USENIX Security 20)*, pages 1605–1622, 2020.

Aleena Anna Thomas, David Ifeoluwa Adelani, Ali Davody, Aditya Mogadala, and Dietrich Klakow. Investigating the Impact of Pre-trained Word Embeddings on Memorization in Neural Networks. In *Workshop on Time-Delay Systems*, 2020. URL `https://api.semanticscholar.org/CorpusID:220658693`.

Kushal Tirumala, Aram Markosyan, Luke Zettlemoyer, and Armen Aghajanyan. Memorization without overfitting: Analyzing the training dynamics of large language models. *Advances in Neural Information Processing Systems*, 35: 38274–38290, 2022.

SK Murakonda and R Shokri. ML Privacy Meter: Aiding Regulatory Compliance by Quantifying the Privacy Risks of Machine Learning., 2007.

Congzheng Song, Thomas Ristenpart, and Vitaly Shmatikov. Machine learning models that remember too much. In *Proceedings of the 2017 ACM SIGSAC Conference on computer and communications security*, pages 587–601, 2017b.

Vitaly Feldman. Does learning require memorization? a short tale about a long tail. In *Proceedings of the 52nd Annual ACM SIGACT Symposium on Theory of Computing*, pages 954–959, 2020.

Nicholas Carlini, Daphne Ippolito, Matthew Jagielski, Katherine Lee, Florian Tramer, and Chiyuan Zhang. Quantifying Memorization Across Neural Language Models. In *The Eleventh International Conference on Learning Representations*, 2023b. URL `https://openreview.net/forum?id=TatRHT_1cK`.

Li Deng. The mnist database of handwritten digit images for machine learning research. *IEEE Signal Processing Magazine*, 29(6):141–142, 2012.

Gal Vardi, Gilad Yehudai, and Ohad Shamir. Gradient methods provably converge to non-robust networks. *Advances in Neural Information Processing Systems*, 35:20921–20932, 2022.

Kaifeng Lyu and Jian Li. Gradient descent maximizes the margin of homogeneous neural networks. *arXiv preprint arXiv:1906.05890*, 2019.

Ziwei Ji and Matus Telgarsky. Directional convergence and alignment in deep learning. *Advances in Neural Information Processing Systems*, 33:17176–17186, 2020.

Daniel Berend and Aryeh Kontorovich. A sharp estimate of the binomial mean absolute deviation with applications. *Statist. Probab. Lett.*, 83(4):1254–1259, 2013. ISSN 0167-7152,1879-2103. doi:10.1016/j.spl.2013.01.023. URL `https://doi.org/10.1016/j.spl.2013.01.023`.

T.N. Bhargava and P.H. Doyle. A geometric study of diversity. *Journal of Theoretical Biology*, 43(2):241–251, 1974. ISSN 0022-5193. doi:https://doi.org/10.1016/S0022-5193(74)80057-3. URL `https://www.sciencedirect.com/science/article/pii/S0022519374800573`.

C.Radhakrishna Rao. Diversity and dissimilarity coefficients: A unified approach. *Theoretical Population Biology*, 21(1):24–43, 1982. ISSN 0040-5809. doi:https://doi.org/10.1016/0040-5809(82)90004-1. URL `https://www.sciencedirect.com/science/article/pii/0040580982900041`.

C.R. Rao, University of Pittsburgh. Institute for Statistics, and Applications. *Gini-Simpson Index of Diversity: A Characterization, Generalization and Applications*. Technical report (University of Pittsburgh. Institute for Statistics and Applications). University of Pittsburgh, 1981. URL `https://books.google.ca/books?id=d7laNQAACAAJ`.

J. Ziv and M. Zakai. On Functionals Satisfying a Data-Processing Theorem. *IEEE Trans. Inf. Theor.*, 19(3):275–283, may 1973. ISSN 0018-9448. doi:10.1109/TIT.1973.1055015. URL `https://doi.org/10.1109/TIT.1973.1055015`.

Vlad Bally and Lucia Caramellino. Asymptotic development for the CLT in total variation distance. *Bernoulli*, 22(4): 2442–2485, 2016. ISSN 1350-7265,1573-9759. doi:10.3150/15-BEJ734. URL `https://doi.org/10.3150/15-BEJ734`.

Luc Devroye, Abbas Mehrabian, and Tommy Reddad. The total variation distance between high-dimensional Gaussians with the same mean. *arXiv preprint arXiv:1810.08693*, 2018.

Abraham De Moivre. *Miscellanea Analytica de Seriebus et Quadraturis*. J. Thonson and J. Watts, London, 1730.

Herbert Robbins. A remark on Stirling's formula. *Amer. Math. Monthly*, 62:26–29, 1955. ISSN 0002-9890,1930-0972. doi:10.2307/2308012. URL `https://doi.org/10.2307/2308012`.

## A   More comments on Section 2

We further discuss here the assumption that the algorithm shall be expressed as a function of the empirical distribution of the training dataset.

Usual definition of an algorithm $\mathcal{A}$ asks for its domain to be $\bigcup_{n \geq 1} \mathcal{Z}^n$ which is similar to identify it with a sequence of algorithms $(\mathcal{A}_n : \mathcal{Z}^n \to \mathcal{P}')_{n \geq 1}$ where for each $n \geq 1$ we have that the restriction of $\mathcal{A}$ to $\mathcal{Z}^n$ coincides with $\mathcal{A}_n$. However, this definition allows not specifying its behaviour through all values of $n$, and specifically having drastically different behaviours for different values of $n$. To rigorously study the characteristics of an algorithm, it is natural to ask that its behaviour is similar for all values of $n$, meaning that is behaviour can be defined independently of $n$.

Our assumption solves this issue, as the function $G : \mathcal{M} \to \mathcal{P}'$ from Section 2 is defined for all discrete distribution on $\mathcal{Z}$. Furthermore, it is worth noting that this assumption holds for all algorithms aiming at minimizing the empirical cost on its training dataset. For most algorithms, it will still hold even when some weights are applied to the samples. Indeed, changing the distribution $P$ by $P'$ for some other distribution $P'$, the training dataset size $n$ by some other integer $n' \in \mathbb{N}$ and adequately adapt the algorithm $\mathcal{A}$ into another algorithm $\mathcal{A}'$ to take into account the changes, makes the study still valid as long as we consider $\Delta_{n'}(P', \mathcal{A}')$ instead of $\Delta_n(P, \mathcal{A})$. It is then sufficient to study under this hypothesis.

This assumption in particular treats all points similarly, and is invariant with respect to the redundancy of the whole dataset. In particular, it is equivalent to saying that the algorithm is **symmetric** and **redundancy invariant**, whose definitions are given below.

**Definition A.1** (Symmetric Map). *Given two sets $\mathcal{Z}_1$ and $\mathcal{Z}_2$ and an integer $k$, a map $f : \mathcal{Z}_1^k \to \mathcal{Z}_2$ is said to be* **symmetric** *if for any $(a_1, \cdots, a_k) \in \mathcal{Z}_1^k$ and any permutation $\sigma$ on $\{1, \cdots, k\}$, we have*

$$f(a_1, \cdots, a_k) = f\left(a_{\sigma(1)}, \cdots, a_{\sigma(k)}\right).$$

**Definition A.2** (Redundancy Invariant Map). *Given two sets $\mathcal{Z}_1$ and $\mathcal{Z}_2$, a map $f : \bigcup_{k>0} \mathcal{Z}_1^k \to \mathcal{Z}_2$ is said to be* **redundancy invariant** *if for any integer $m$ and any $\boldsymbol{a} = (a_1, \cdots, a_m) \in \mathcal{Z}_1^m$, we have*

$$f(\boldsymbol{a}) = f(\boldsymbol{a}, \cdots, \boldsymbol{a}).$$

We summarize the last claim in the following proposition.

**Proposition A.3.** *Let $f : \bigcup_{k>0} \mathcal{Z}^k \to \mathcal{Z}'$ be a measurable map onto any space $\mathcal{Z}'$. Then the following statements are equivalent*

> *(i) $f$ is redundancy invariant and for any $k \in \mathbb{N}$, the restriction of $f$ to $\mathcal{Z}^k$ is symmetric.*

> *(ii) There exists a function $G : \mathcal{M} \to \mathcal{Z}'$ such that for any $k \in \mathbb{N}$, for any $(z_1, \cdots, z_k) \in \mathcal{Z}^k$ we have $f(z_1, \cdots, z_k) = G\left(\frac{1}{k} \sum_{j=1}^k \delta_{z_j}\right)$.*

*Proof of Proposition A.3.* We only prove that $(i)$ implies $(ii)$. The fact that $(ii)$ implies $(i)$ is straightforward.

Let $f : \bigcup_{k>0} \mathcal{Z}^k \to \mathcal{Z}'$ be a measurable map satisfying condition $(i)$. Let $\mathcal{M}^{\text{emp}}$ be the set of all possible empirical distributions, that is the subset of $\mathcal{M}$ containing all $\frac{1}{k} \sum_{j=1}^k \delta_{z_j}$ for all integer $k$ and all $(z_1, \cdots, z_k) \in \mathcal{Z}^k$. We shall define $G$ on $\mathcal{M}^{\text{emp}}$ such that $(ii)$ holds true.

For any $Q \in \mathcal{M}^{\text{emp}}$, let $\{z_1, \cdots, z_m\}$ be its support and $q_1, \ldots, q_m \in (0, 1)$ be such that $Q = \sum_{j=1}^m q_j \delta_{z_j}$. Since $Q$ is an empirical distribution, there exists positive integers $k_1, \ldots, k_m$ (for each $j$, $k_j$ is the number of occurences of $z_j$ in the sample from which $Q$ is the empirical distribution) such that $q_j = \frac{k_j}{K}$, with $K = \sum_{j=1}^m k_j$.

Let $r = gcd(k_1, \ldots, k_m)$ be the greatest common divisor of the $k_j$'s and define $k'_j = k_j/r$ for $j = 1, \ldots, m$. Then with $K' := \sum_{j=1}^m k'_j$, we have $q_j = \frac{k'_j}{K'}$.

Now, for any other sequence of positive integers $\ell_1, \ldots, \ell_m$ such that $q_j = \frac{\ell_j}{L}$, with $L = \sum_{j=1}^m \ell_j$, we get for all $j$, $\ell_j = sk'_j$ with $s = gcd(\ell_1, \ldots, \ell_m)$. Thus we may define $G(Q) = f(\boldsymbol{z})$ where $\boldsymbol{z}$ is the dataset consisting of all $z_j$'s with $k'_j$ repetitions.

We now prove that such a $G$ satisfies $(ii)$. Indeed, for any integer $k$ and any $Z := (z'_1, \cdots, z'_k) \in \mathcal{Z}^k$, define $V := ((\ell_1, z_1), \cdots, (\ell_m, z_m))$ where $(z_1, \cdots, z_m)$ are the distinct elements of $Z$ and $(\ell_1, \cdots, \ell_m)$ are their occurrences. Define $r$ as their greatest common divisor, and $(k_1, \ldots, k_m) = (\ell_1, \cdots, \ell_m)/r$. By using the fact that $f$ is symmetric and redundancy invariant, we get that $f(Z) = f(Z_0) = G(Q)$ where $Z_0$ is the dataset consisting of all $z_j$'s with $k_j$ repetitions and $Q = \sum_{j=1}^m \frac{k_j}{K} \delta_{z_j} = \frac{1}{n} \sum_{j=1}^n \delta_{z'_j}$. Thus $(ii)$ holds. $\square$

# B    More comments on Overfitting

We give here more details about Section 4. Specifically, we further discuss Proposition 4.2, we give an extension of Theorem 4.3 to the setting of classification.

## B.1    Additional comments on Proposition 4.2 and Theorem 4.3.

Proposition 4.2 together with the second point of Theorem 4.3 state that if we train our regressor long enough so that with high probability, the training loss of one training sample is below $\varepsilon$ for some tiny threshold $\varepsilon$, then the probability that the loss of an independent sample reaches this threshold is near 0. In particular it confirms an already well-known empirical insight that overfitting implies poor security.

Interestingly, if we only assume Definition 4.1 to hold without Proposition 4.2 to hold, then a much weaker version of the second point of Theorem 4.3 still holds. Indeed, for a fixed $\alpha \in (0,1)$, given a sequence of algorithms $(\mathcal{A}^\varepsilon)_{\varepsilon \in \mathbb{R}^+}$ that are $(\varepsilon, 1-\alpha)$-overfitting for all $\varepsilon > 0$, we have that $\lim_{\varepsilon \to 0} \Delta_n(P, \mathcal{A}^\varepsilon) \geq 1 - \alpha$.

## B.2    Extension to classification

The second point of Theorem 4.3 requires the absolute continuity of the distribution of the label with respect to the Lebesgue measure, which makes it not straightforward to extend it to classifiers.
We discuss here one very specific framework in which we have been able to extend our results to the classification setting. The framework and the assumptions are all inspired from Vardi et al. [2022].

We assume that the data space is restrained to the binary classification setting with data in the sphere of radius $\sqrt{s}$, i.e. $\mathcal{Z} := \left(\sqrt{s}\mathbb{S}^{s-1}\right) \times \{-1, 1\}$ where $\mathbb{S}^{s-1}$ is the unit sphere in $\mathbb{R}^s$. We assume our data $(z_1, \cdots, z_n) :=$ $((x_1, y_1), \cdots, (x_n, y_n))$ to be independently drawn on $\mathcal{Z}$ from a distribution $P$. We assume that the conditional law of $x_1$ given $y_1$ is absolutely continuous with respect to the Lebesgue measure on $\sqrt{s}\mathbb{S}^{s-1}$. We denote by $\mathcal{H}$ the latter hypothesis. Let $\Psi_\theta(x) = \sum_{j=1}^l v_j \sigma(w_j^T x + b_j)$ be a 2−ReLU network with parameters $\theta$, i.e. $\theta = (v_j, w_j, b_j)_{j=1}^l$ with $l \in \mathbb{N}$ the width of the network and $\sigma(u) = \max(u, 0)$. We aim at learning a classifier $\Psi_{\hat{\theta}_n}$ on the data by minimizing

$$\mathcal{L} : \theta \mapsto \sum_{j=1}^n l(y_j \Psi_\theta(x_j)), \tag{12}$$

where $l : \mathbb{R} \to \mathbb{R}^+$ is either the exponential loss or the logistic loss. To reach the objective, we apply Gradient Flow on the objective Equation 12, producing a trajectory $\theta_n(t)$ at time $t$. From Vardi et al. [2022] Theorem 3.1, there exists a 2−ReLU network classifying perfectly the training dataset, as long as $\max_{i \neq j}\left\{|x_i^T x_j|\right\} < d$, which holds almost surely by $\mathcal{H}$. Let the initial point $\theta_n(0)$ be the parameters of this network.

Then by Vardi et al. [2022] Theorem 2.1, paraphrasing Lyu and Li [2019], Ji and Telgarsky [2020], $\frac{\theta_n(t)}{\|\theta_n(t)\|}$ converges as $t$ tends to infinity to some vector $\bar{\theta}_n$ which is colinear to some KKT point of the following problem

$$\min_\theta \frac{1}{2}\|\theta\|^2 \text{ s.t. } \forall i = 1, \ldots, n; y_i \Psi_\theta(x_i) \geq 1. \tag{13}$$

Conditional to the event $E := \text{"}\max_{i \neq j}\left\{|x_i^T x_j|\right\} \leq \frac{s+1}{3n} - 1\text{"}$, by Vardi et al. [2022] Lemma C.1 we get that for all $j = 1, \ldots, n$, we have

$$y_j \Psi_{\bar{\theta}_n}(x_j) = \lambda(z_1, \cdots, z_n), \tag{14}$$

for some $\lambda(z_1, \cdots, z_n) > 0$.

We consider our algorithm $\mathcal{A}$ to output

$$\mathcal{A}(z_1, \cdots, z_n) = \hat{\theta}_n := \frac{\bar{\theta}_n}{\sqrt{\lambda(z_1, \cdots, z_n)}},$$

which gives the same classifier as with $\bar{\theta}_n$.

We then get the following result.

**Proposition B.1.** *Assume that $l \geq n$ and let $C := \max_{i \neq j} \{|x_i^T x_j|\}$. Then, there exists an initialization $\theta_n(0)$ of the gradient flow for which it holds that*

$$\Delta_n(P, \mathcal{A}) \geq P\left(C \leq \frac{s+1}{3n} - 1\right).$$

*Moreover, if the marginal distribution of $x$ is the uniform distribution on $\sqrt{s}\mathbb{S}^{s-1}$, then*

$$\Delta_n(P, \mathcal{A}) \geq 1 - s^{3 - ln(s)/4},$$

*as soon as $n \leq \frac{1}{3} \frac{s+1}{\sqrt{s}ln(s)+1}$.*

*Proof of Proposition B.1.* By definition of $\Psi_\theta$ for any $\theta \in \Theta$, it holds that these networks are $2-$homogeneous, so that conditional to the event $E$, Equation 14 leads to

$$y_j \Psi_{\hat{\theta}_n}(x_j) = 1, \tag{15}$$

for any $j = 1, \ldots, n$.

Let $S := \{(\theta, x, y) \in \Theta \times \left(\sqrt{s}\mathbb{S}^{s-1}\right) \times \{-1, 1\} : y\Psi_\theta(x) = 1\}$. Then, by definition of $\Delta_n(P, \mathcal{A})$, we have

$$\begin{aligned}
\Delta_n(P, \mathcal{A}) &\geq P((\hat{\theta}_n, x_1, y_1) \in S) - P((\hat{\theta}_n, x, y) \in S) \\
&= P((\hat{\theta}_n, x_1, y_1) \in S \mid E)P(E) + P((\hat{\theta}_n, x_1, y_1) \in S \mid E^c)P(E^c) - \mathbb{E}\left[P(\Psi_{\hat{\theta}_n}(x) = y \mid \hat{\theta}_n, y)\right] \\
&\geq P((\hat{\theta}_n, x_1, y_1) \in S \mid E)P(E) - \mathbb{E}\left[P(\Psi_{\hat{\theta}_n}(x) = y \mid \hat{\theta}_n, y)\right],
\end{aligned}$$

where we have lower bounded the second term by 0.

By Equation 15, we have $P((\hat{\theta}_n, x_1, y_1) \in S \mid E) = 1$. Now, by independence between $(x, y)$ and $\hat{\theta}_n$, it is sufficient to show that for any $\theta \in \Theta$, we have $P(\Psi_\theta(x) = y \mid y) = 0$ almost surely. Without loss of generality, we may assume that $v_j \neq 0$ for any $j = 1, \ldots, l$. We set $B_J(x, y) := \{\forall j \in J, w_j^T x + b_j > 0\} \cap \{\forall j \in J^c, w_j^T x + b_j \leq 0\} \cap \left\{\sum_{j \in J} v_j\left(w_j^T x + b_j\right) = y\right\}$ for any $J \subseteq [1, \cdots, l]$. We then get

$$P\left(\Psi_{\hat{\theta}_n}(x) = y \mid y\right) = \sum_{J \subseteq [1, \cdots, l]} P\left(B_J(x, y) \mid y\right)$$

$$\leq \sum_{J \subseteq [1, \cdots, l]} P\left(\sum_{j \in J} v_j\left(w_j^T x + b_j\right) = y \mid y\right).$$

Note that the space $H_{y,J} := \left\{x \in \mathbb{R}^s : \sum_{j \in J} v_j\left(w_j^T x + b_j\right) = y\right\}$ is an hyperplan of $\mathbb{R}^s$ for any $y \in \{-1, 1\}$ and any $J \subseteq [1, \cdots, l]$. Then the quantity $P(x \in H_{y,J} \mid y)$ equals 0 by $\mathcal{H}$. Hence,

$$\Delta_n(P, \mathcal{A}) \geq P(E).$$

Under the further hypothesis that $x$ is uniformly distributed on the sphere, and that $n \leq \frac{1}{3} \frac{s+1}{\sqrt{s}ln(s)+1}$, it holds that $\frac{s+1}{3n} - 1 \geq \frac{\sqrt{s}}{ln(s)}$. Then Vardi et al. [2022] Lemma 3.1 concludes.

$\square$

## C   More comments on Section 5

We give here some more details about the behaviour of $\Delta_n(P, \mathcal{A})$ when the set of parameters $\Theta$ has finite cardinal. We also further discuss the quantity $C(P)$.

### C.1   Different rates for $\Delta_n(P, \mathcal{A})$

Corollary 5.4.1 gives a rate of $n^{-1/2}$ for $\Delta_n(P, \mathcal{A})$ when $C(P) < \infty$ and $n$ is sufficiently large. In the case when $C(P)$ is infinite, it is interesting to note that we still have convergence to 0 of $\Delta_n(P, \mathcal{A})$ but at an arbitrarily slow rate. We formalize this result in the following lemma :

**Lemma C.1.** *If* $C(P) = \infty$, $\max_{\mathcal{A}} \Delta_n(P, \mathcal{A})$ *still tends to* 0 *as* $n$ *tends to infinity, but the (depending on P) rate can be arbitrarily slow.*

In this case, in order to find the minimum amount of data to get a control on $Sec_n(P, \mathcal{A})$ requires the estimation of the r.h.s. of Equation 11 which is not obvious, as the condition $C(P) = \infty$ is equivalent to $K = \infty$.

*Proof.* It is a direct corollary of Lemmas 7 and 8 of Berend and Kontorovich [2013]. $\qquad\square$

Further observe that Theorem 5.4 is valid only when the support $\Theta$ of the parameters has infinite cardinal. Indeed, as demonstrated in the proof of Theorem 5.4 in Section G, the result holds if there exist injective maps from $\mathcal{Z}$ into $\Theta$. In the case when $\Theta$ has finite cardinal $L \in \mathbb{N}$, it is hard to get insightful formulas similar to Equation 11, howerver it is possible to rewrite it as follows

**Lemma C.2.** *Let* $\mathbf{r} := (N_1, \cdots, N_K)$ *be a random vector having multinomial distribution with parameters* $(n; p_1, \cdots, p_K)$. *There exists a partition* $(D_l)_{l=1\ldots L}$ *of the support of* $\mathbf{r}$ *such that*

$$\Delta_n(P, \mathcal{A}) = \frac{1}{2} \sum_{j=1}^{K} \sum_{l=1}^{L} \left| \mathbb{E}\left[ \left\{ p_j - \frac{N_j}{n} \right\} 1\{\mathbf{r} \in D_l\} \right] \right|.$$

Lemma C.2 is a tool to understand the behaviour of $\Delta_n(P, \mathcal{A})$ depending on the structure of the algorithm $\mathcal{A}$. Although there is a strong similarity between Lemma C.2 and Theorem 5.4, the value of $\Delta_n(P, \mathcal{A})$ in Lemma C.2 is smaller than the right hand side of Equation 11. This could informally mean that discretizing/quantizing an algorithm improves its security.

In the case of finite parameters space $\Theta$, it is easy to come up with an example for which the rate of $\Delta_n(P, \mathcal{A})$ is exponential.

**Lemma C.3.** *Let* $P$ *be the Bernoulli distribution with parameter* $p \in (0, 1)$ *and let* $\hat{\theta}_n := \sup_j z_j$. *Then,*

$$\Delta_n(P, \mathcal{A}) = 2p(1 - p)^n.$$

*Proof of Lemma C.2.* We recall that the range of the algorithm is finite. Without loss of generality, we identify the set of parameters to $\{1, \cdots, L\}$ for some $L \in \mathbb{N}$. Here Equation 29 writes

$$2\Delta_n(P, \mathcal{A}) = \sum_{k=1}^{K} p_k \sum_{l=1}^{L} \left| P(F(\hat{P}_n) = l) - P(F(\hat{P}_n^k) = l) \right|. \tag{16}$$

By the symmetry of $\mathcal{A}$, only the distribution of the data among the possible values is relevant, that is the random variables $N_1 := \sum_{i=1}^{n} \mathbf{1}_{z_i=1}, \cdots, N_K := \sum_{i=1}^{n} \mathbf{1}_{z_i=K}$. Note that for $j = 1, \cdots, K$, the random variable $N_j$ is the number of occurrences of $u_j$ in the dataset.
By the *i.i.d.* assumption of the data, the random vector $\mathbf{r} := (N_1, \cdots, N_K)$ follows a multinomial distribution with parameters $(n; p_1, \cdots, p_K)$.
For any $k = 1, \cdots, K$, let $D_k \subseteq \{(n_1, \cdots, n_K) \in \{0, \cdots, n\}^K : \sum_{q=1}^{K} n_q = n\}$ such that $F(\hat{P}_n) = k$ if and only if $\mathbf{r} \in D_k$.
Since $\mathbf{r}$ follows a multinomial distribution, using Equations 16, 31 and 32, one gets

$$2\Delta_n(P,\mathcal{A}) = \sum_{k=1}^{K} p_k \sum_{l=1}^{L} |P(\mathbf{r} \in D_l) - P(\mathbf{r} \in D_l \mid \mathbf{z}_1 = u_k)|$$

$$= \sum_{k=1}^{K} p_k \sum_{l=1}^{L} \left| \sum_{(n_1,\cdots,n_K) \in D_l} \binom{n}{n_1 \cdots n_K} \left( \prod_{q=1}^{K} p_q^{n_q} \right) \left\{ 1 - \frac{n_k}{np_k} \right\} \right|$$

$$= \sum_{k=1}^{K} \sum_{l=1}^{L} \left| \mathbb{E} \left[ \left\{ p_k - \frac{N_k}{n} \right\} \mathbf{1}_{\mathbf{r} \in D_l} \right] \right|.$$

$\square$

*Proof of Lemma C.3.* From Equation 24, one has

$$\Delta_n(P,\mathcal{A}) = \mathbb{E} \left( \left\| \mathcal{L}(F(\hat{P}_n)) - \mathcal{L}(F(\hat{P}_n)|\mathbf{z}_1) \right\|_{\mathrm{TV}} \right).$$

By definition of the total variation for discrete distributions, for $b \in \{0,1\}$, one has

$$2 \left\| \mathcal{L}(F(\hat{P}_n)) - \mathcal{L}(F(\hat{P}_n)|\mathbf{z}_1 = b) \right\|_{\mathrm{TV}} = |P(F(\hat{P}_n) = 1) - P(F(\hat{P}_n) = 1|\mathbf{z}_1 = b)|$$

$$+ |P(F(\hat{P}_n) = 0) - P(F(\hat{P}_n) = 0|\mathbf{z}_1 = b)|$$

$$= \left| 1 - (1-p)^n - \left\{ \begin{array}{ll} 1 & \text{si } b = 1 \\ 1 - (1-p)^{n-1} & \text{si } b = 0 \end{array} \right. \right|$$

$$+ \left| (1-p)^n - \left\{ \begin{array}{ll} 0 & \text{if } b = 1 \\ (1-p)^{n-1} & \text{if } b = 0 \end{array} \right. \right|$$

$$= 2(1-p)^{n-1} |(1-p) - \mathbf{1}_{b=0}|$$

$$= 2(1-p)^{n-1} |p - \mathbf{1}_{b=1}|.$$

Taking expectation over $\mathbf{z}_1$ gives

$$\Delta_n(P,\mathcal{A}) = (1-p)^{n-1} \mathbb{E}[|p - \mathbf{1}_{\mathbf{z}_1=1}|]$$

$$= (1-p)^{n-1} [2p(1-p)]$$

$$= 2p(1-p)^n,$$

which concludes the proof. $\square$

## C.2 Relation of $C(P)$ with the Gini-Simpson Entropy and the Shannon's Entropy

We recall that for a discrete random variable $X$ with distribution $P = \sum_{j=1}^{K} p_j \delta_{u_j}$, the Gini-Simpson Entropy is given by (see Bhargava and Doyle [1974], Rao [1982])

$$\text{G-S}(X) := 1 - \sum_{i=1}^{K} p_i^2,$$

and the Shannon Entropy is given by

$$H(X) := - \sum_{i=1}^{K} p_i \log(p_i).$$

From the inequality $\frac{1}{2}\sqrt{p(1-p)} \geq p(1-p)$ for all $p \in [0,1]$, we have

$$\frac{C(P)}{2} = \frac{1}{2} \sum_{i=1}^{K} \sqrt{p_i(1-p_i)} \geq \sum_{i=1}^{K} p_i(1-p_i) = \text{G-S}(X). \tag{17}$$

From the concavity of the square root, we also have

$$\frac{C(P)}{K} = \frac{1}{K} \sum_{i=1}^{K} \sqrt{p_i(1-p_i)} \leq \sqrt{\frac{1}{K} \sum_{i=1}^{K} p_i(1-p_i)} = \sqrt{\frac{1}{K} \text{G-S}(X)}. \tag{18}$$

The Gini-Simpson index can be interpreted as the expected distance between two randomly selected individuals when the distance is defined as zero if they belong to the same category and one otherwise [Rao et al., 1981], that is $\mathbb{P}(X \neq Y)$ for $X$ and $Y$ i.i.d.. The inequality mentioned above suggests that as the Gini-Simpson index increases (e.g., higher diversity of the data), security decreases and thus, the MIAs are expected to be more successful. Another, such commonly used diversity measure is Shannon entropy. Interestingly, $C(P)$ can be also upped and lower bounded by the Shannon entropy as follows:

$$H(X) \leq C(P) \leq \sqrt{K} \sqrt{H(X)}. \tag{19}$$

These bounds easily follow by noticing that

$$C(P) \leq \sqrt{K \left[1 - \exp\left(-H(X)\right)\right]} \leq \sqrt{KH(X)},$$

by upper bounding:

$$-\log\left(\sum_{i=1}^{K} p_i^2\right) \leq H(X).$$

Similarly,

$$-p_i \log p_i \leq \sqrt{p_i(1-p_i)}, \quad \text{for all } 0 \leq p_i \leq 1$$

and thus, $C(P) \geq H(X)$.

The Gini-Simpson Entropy and the Shannon Entropy are maximized by the uniform distribution. This is also the case for $C(P)$, as proved below.

**Lemma C.4.** *Let $P := \sum_{j=1}^{K} p_j \delta_{u_j}$ be a discrete distribution with finite $K$. Let $\mathcal{M}_K$ be the set of all such distributions. We then have the following properties on $C(P)$. For fixed $K \geq 2$, we have*

- $\max\limits_{P \in \mathcal{M}_K} C(P) = \sqrt{K-1}$

- $\underset{P \in \mathcal{M}_K}{argmax} \, C(P) = Unif(u_1, \cdots, u_K).$

Interestingly, one can observe that for a fixed and finite number of atoms $K$, the sub-levels of $C(P)$ are tightly controlled by the value of $\max\limits_{j} p_j$. More precisely, at fixed $\max\limits_{j} p_j := \delta$ for some fixed $\delta \in [1/K; 1)$, the width of the interval $(inf\{C(P)\}, max\{C(P)\}]$ is entirely determined by $\delta$ and $K$. Specifically, the further $\delta$ is from $1/2$, the thinner the interval gets. We summarize this comment below.

**Lemma C.5.** *Let $\delta \in [1/K, 1)$. Then the following statements hold :*

$$\max\left\{C(P) : P \in \mathcal{M}_K, \max_j p_j = \delta\right\} = \sqrt{\delta(1-\delta)} + \sqrt{(1-\delta)(K-2+\delta)}$$

$$\inf\left\{C(P) : P \in \mathcal{M}_K, \max_j p_j = \delta\right\} = \sqrt{\delta(1-\delta)}\lfloor \delta^{-1} \rfloor + \sqrt{\delta\lfloor \delta^{-1} \rfloor(1 - \delta\lfloor \delta^{-1} \rfloor)}$$

*Proof of Lemma C.5.* Let $\mathcal{M}_K(\delta) := \{P \in \mathcal{M}_K : \max_j p_j = \delta\}$. Without loss of generality, we always assume that $p_K = \max_j p_j$. Let $f : p \mapsto \sqrt{p(1-p)}$.

First notice that $f$ is a concave function, so that for any $p_1, \cdots, p_m \in [0,1]$ we have $\frac{1}{m} \sum_{j=1}^{m} f(p_j) \leq f\left(\frac{1}{m} \sum_{j=1}^{m} p_j\right)$.

In particular, for any $P \in \mathcal{M}_K(\delta)$, we have $\sum_{j=1}^{K-1} p_j = 1 - \delta$ giving

$$C(P) \leq C\left(\underbrace{\frac{1-\delta}{K-1}, \cdots, \frac{1-\delta}{K-1}}_{K-1}, \delta\right).$$

17

Evaluating the r.h.s. of the last inequality gives us the first result.

If $\delta \geq 1/2$, using again the concavity of $f$ gives us that for any $P \in \mathcal{M}_K(\delta)$, we have

$$C(P) \geq C(\underbrace{0, \cdots, 0}_{K-2}, 1 - \delta, \delta),$$

where the last quantity evaluates to $2\sqrt{\delta(1-\delta)}$.
If $\delta < 1/2$, using the concavity of $f$ gives us that for any $P \in \mathcal{M}_K(\delta)$, we have

$$C(P) \geq C(0, \cdots, 0, 1 - \delta\lfloor\delta^{-1}\rfloor, \underbrace{\delta, \cdots, \delta}_{\lfloor\delta^{-1}\rfloor}),$$

evaluating at $\sqrt{\delta(1-\delta)}\lfloor\delta^{-1}\rfloor + \sqrt{\delta\lfloor\delta^{-1}\rfloor(1 - \delta\lfloor\delta^{-1}\rfloor)}$ for any numbers of 0s. Combining the two results give the last equality of the lemma. $\qquad\square$

*Proof of Lemma C.4.* Denote by $M_K(\delta) := \sqrt{\delta(1-\delta)} + \sqrt{(1-\delta)(K-2+\delta)}$. By Lemma C.5, we have that

$$\max_{P \in \mathcal{M}_K} C(P) = \max_{\delta \in [1/K, 1)} M_K(\delta) = \sqrt{K-1} \text{ reached at } \Big(\underbrace{\frac{1}{K}, \cdots, \frac{1}{K}}_{K}\Big) \in \mathcal{M}_K(1/K). \qquad\square$$

# D  Setting for Section 6

We provide here additional values used for the numerical experiments presented in Section 6. All experiments have been conducted with PyTorch library.

## D.1  Details for Section 6.1

We consider here synthetically generated data $x_1, \cdots, x_n \overset{i.i.d.}{\sim} \mathrm{Unif}(\mathbb{S}_{d-1})$ with $d = 10$ and $n = 1000$. For any $i = 1, \ldots, n$, we set $y_i = \Psi^*(x_i) + \zeta_i$, where $\zeta_i \sim \mathcal{N}(0, 0.01)$ is some independent noise and $\Psi^*$ is defined for all $x \in \mathbb{R}^d$ by $\Psi^*(x) = \sin(\pi \beta^T x)$ for some fixed $\beta \in \mathbb{S}_{d-1}$. For the sake of simplicity, we took $\beta = (1, 0, \cdots, 0)$, as the distribution of the data is rotation-invariant. We aim at estimating $\Psi^*$ by a 2 layers ReLU neural network, whose hidden layer has width 4096. We train the neural network $\Psi_\theta$ by minimizing the MSE loss with the *Adam* optimizer and learning rate 0.1 for 2500 iterations.

To evaluate the accuracy, we generated 10000 test samples $x_1^{test}, \cdots, x_{10000}^{test} \overset{i.i.d.}{\sim} \mathrm{Unif}(\mathbb{S}_{d-1})$ independently from the training dataset.

## D.2  Details for Section 6.2

We consider here the whole MNIST dataset, and use the given separation between training and test. We aimed at classifying the dataset by learning on a 3-layers ReLU neural network, with both internal layers having width 256. We trained the neural networks by minimizing the cross-entropy loss with the *Adam* optimizer and learning 0.01 for 500 iterations. To create the discretizations, we drew three times 1000 samples from the dataset (that have been used only for that purpose) and we constructed the clusterizing algorithm based on the **MiniBatchKMeans** function from *scikit-learn* library. We used $K = 20, 50, 100$ clusters respectively. For each dataset size $n = 1000, 5000, 10000$, we performed the following steps :

- We draw a dataset $D_n$ of size $n$ not containing the samples used for the clusterings.
- We trained a first neural network on $D_n$ (column **raw dataset**)
- For each clustering, we discretized $D_n$ and then trained a neural network on the new dataset.

# E  Proofs of Section 3

We give here the proof of Theorem 3.1.

*Proof of Theorem 3.1.* From the law of total probability, we have

$$\text{Acc}_n(\phi; P, \mathcal{A}) = P(\phi(\hat{\theta}_n, \tilde{z}) = 1 - T)$$

$$= P(T = 1)P\left(\phi(\hat{\theta}_n, \tilde{z}) = 1 - T | T = 1\right) + P(T = 0)P\left(\phi(\hat{\theta}_n, \tilde{z}) = 1 - T | T = 0\right)$$

$$= \nu P\left(\phi(\hat{\theta}_n, z_0) = 0\right) + (1 - \nu)P\left(\phi(\hat{\theta}_n, z_1) = 1\right),$$

where the third equality comes from the definition of $\tilde{z}$ and $T$. We now define $B := \{(\theta, z) \in \Theta \times \mathcal{Z} : \phi(\theta, z) = 1\}$ and rewrite $\text{Acc}_n(\phi; P, \mathcal{A})$ as

$$\text{Acc}_n(\phi; P, \mathcal{A}) = \nu \left(1 - P\left((\hat{\theta}_n, z_0) \in B\right)\right) + (1 - \nu)P\left((\hat{\theta}_n, z_1) \in B\right). \tag{20}$$

Taking the maximum over all MIAs $\phi$ then reduces to taking the maximum of the r.h.s. of Equation 20 over all measurable sets $B$. Setting $\gamma := \frac{\nu}{1 - \nu}$, we then get

$$\max_{\phi} \text{Acc}_n(\phi; P, \mathcal{A}) = (1 - \nu)\max_{B}\left[P\left((\hat{\theta}_n, z_0) \in B\right) - \gamma P\left((\hat{\theta}_n, z_1) \in B\right)\right] + \nu, \tag{21}$$

where the maximum is taken over all measurable sets $B$. Let now $\zeta$ be a dominating measure of the distributions of $(\hat{\theta}_n, z_0)$ and $(\hat{\theta}_n, z_1)$ (for instance their average). We denote by $p$ (resp. $q$) the density of the distribution of $(\hat{\theta}_n, z_0)$ (resp. $(\hat{\theta}_n, z_1)$) with respect to $\zeta$. Then, the involved maximum in the r.h.s. of Equation 21 is reached on the set

$$B^* := \{p/q \geq \gamma\}.$$

The maximum being taken over all measurable sets in Equation 21, we may consider replacing $B$ by its complementary $B^c$ in the expression giving

$$\max_{\phi} \text{Acc}_n(\phi; P, \mathcal{A}) = (1 - \nu)\max_{B}\left[\gamma P\left((\hat{\theta}_n, z_1) \in B\right) - P\left((\hat{\theta}_n, z_0) \in B\right)\right] + (1 - \nu), \tag{22}$$

where in this case the maximum is reached on the set

$$B^{*c} := \{p/q < \gamma\}.$$

Taking the average on Equations 21 and 22, we get

$$\max_{\phi} \text{Acc}_n(\phi; P, \mathcal{A}) = \frac{1}{2} + \frac{1}{2}\int \left|(1 - \nu)p - \nu q\right| d\zeta. \tag{23}$$

By the triangular inequality, we may obtain the two following inequalities:

$$\max_{\phi} \text{Acc}_n(\phi; P, \mathcal{A}) \leq \frac{1}{2} + \frac{|1 - 2\nu|}{2}\int q\, d\zeta + \frac{1 - \nu}{2}\int |p - q|\, d\zeta,$$

$$\max_{\phi} \text{Acc}_n(\phi; P, \mathcal{A}) \leq \frac{1}{2} + \frac{|1 - 2\nu|}{2}\int p\, d\zeta + \frac{\nu}{2}\int |p - q|\, d\zeta.$$

With $\int q\, d\zeta = \int p\, d\zeta = 1$, it holds that when $\nu \leq 1/2$, we have $1 - 2\nu \geq 0$ so that $1/2 + |1 - 2\nu|/2 = 1 - \nu$. Similarly, we get $1/2 + |1 - 2\nu|/2 = \nu$ when $\nu \geq 1/2$. Then, setting $\nu_* := \min\{\nu, 1 - \nu\}$ we have in both cases

$$1/2 + |1 - 2\nu|/2 = 1 - \nu_*.$$

Since $\Delta_n(P, \mathcal{A}) = (1/2)\int |p - q|\, d\zeta$ from the definition of the total variation distance, by taking the minimum over the two previous expressions, we have

$$\max_{\phi} \text{Acc}_n(\phi; P, \mathcal{A}) \leq 1 - \nu_* + \nu_* \Delta_n(P, \mathcal{A}),$$

from which we deduce

$$\text{Sec}_n(P, \mathcal{A}) \geq 1 - \Delta_n(P, \mathcal{A}).$$

Following the same steps for the minimum, we have

$$\min_{\phi} \text{Acc}_n(\phi; P, \mathcal{A}) \geq \nu_* - \nu_* \Delta_n(P, \mathcal{A}),$$

hence Theorem 3.1. Equation 23 gives the equality case by plugging $\nu = 1/2$ into it. $\qquad\square$

## F  Proofs of Section 4

We give here the proofs for the Section 4.

*Proof of Proposition 4.2.* Let $l_i := l_i(\hat{\theta}_n) = l_{\hat{\theta}_n}(x_i, y_i)$. The algorithm $\mathcal{A}_{\varepsilon,\alpha}$ stops as soon as $\frac{1}{n} \sum_{j=1}^{n} l_{\hat{\theta}_n}(x_i, y_i) \leq \varepsilon\alpha$.

Let $B_\varepsilon := \left\{ j : l_{\hat{\theta}_n}(x_i, y_i) \leq \varepsilon \right\}$ be the set of samples with loss not larger than $\varepsilon$ at the end of the training. We then have the following sequence of inequalities:

$$n\varepsilon\alpha \geq \sum_{j=1}^{n} l_{\hat{\theta}_n}(x_i, y_i) \geq \sum_{j \in B_\varepsilon} l_{\hat{\theta}_n}(x_i, y_i) + (n - \#B_\varepsilon)\varepsilon \geq (n - \#B_\varepsilon)\varepsilon.$$

From the two extremes, we get that $\sum_{j=1}^{n} 1\{l_j \leq \varepsilon\} := \#B_\varepsilon \geq n(1-\alpha)$. From the *i.i.d.* hypothesis on the data $z_1, \cdots, z_n$, taking the expectation gives the result.

$\square$

*Proof of Theorem 4.3.* We begin by proving the first point. Let $\mathcal{A}$ be an $(\varepsilon, 1-\alpha)$-overfitting algorithm, and let $S^\varepsilon := \{(\theta, x, y) : l_\theta(x, y) \leq \varepsilon\}$. From the definition of $\Delta_n(P, \mathcal{A})$, we have that

$$\begin{aligned}
\Delta_n(P, \mathcal{A}) &\geq P((\hat{\theta}_n, x_1, y_1) \in S^\varepsilon) - P((\hat{\theta}_n, x, y) \in S^\varepsilon) \\
&= P((x_1, y_1) \in S_{\hat{\theta}_n}^\varepsilon) - P((x, y) \in S_{\hat{\theta}_n}^\varepsilon) \\
&= 1 - \alpha - P((x, y) \in S_{\hat{\theta}_n}^\varepsilon) \\
&= 1 - \alpha - \int_{\theta \in \Theta} P((x, y) \in S_\theta^\varepsilon) d\mu_{\hat{\theta}_n},
\end{aligned}$$

which proves the first point.

Now assume that we have a sequence of algorithms $(\mathcal{A}_\eta)_{\eta \in \mathbb{R}^+}$ that stop as soon as $L_n \leq \eta$. Assume the additional hypotheses given in the second point of Theorem 4.3 hold. Let $\alpha \in (0, 1)$ be a fixed scalar. By Proposition 4.2, $\mathcal{A}_\eta$ is $(\eta/\alpha, 1-\alpha)$-overfitting, so that by the first point proven above, we have

$$\Delta_n(P, \mathcal{A}_\eta) \geq 1 - \alpha - \int_{\theta \in \Theta} P((x, y) \in S_\theta^{\eta/\alpha}) d\mu_{\hat{\theta}_n}.$$

For any $\theta \in \Theta$, we have

$$\begin{aligned}
P((x, y) \in S_\theta^{\eta/\alpha}) &= \mathbb{E}\left[ P((x, y) \in S_\theta^{\eta/\alpha} \mid x) \right] \\
&= \mathbb{E}\left[ P(\omega(y, \Psi_\theta(x)) \leq \eta/\alpha \mid x) \right] \\
&\stackrel{\eta \to 0^+}{\to} 0,
\end{aligned}$$

where the limit comes from the continuity of $\omega$ and the absolute continuity of the distribution of y given x. In particular, for any $\alpha \in (0, 1)$, we then have

$$\lim_{\eta \to 0} \Delta_n(P, \mathcal{A}_\eta) \geq 1 - \alpha.$$

Taking the supremum over $\alpha$ gives the result.

$\square$

## G  Proofs of Section 5

*Proof of Theorem 5.1.* Let $m_j := \mathbb{E}\left[ \left\| C^{-1/2} \left\{ L(z_1) - \mathbb{E}[L(z_1)] \right\} \right\|_2^j \right]$ for any positive integer $j$, where $C$ is the covariance matrix of $L(z_1)$. that is the expectation of the $j$-th power of the norm of the centered and reduced version of $L(z_1)$, and $C$ be the covariance matrix of $L(z_1)$.

Setting $L_n := \frac{1}{n} \sum_{j=1}^n L(\mathbf{z}_j)$, by the data processing inequality [Ziv and Zakai, 1973] applied to the total variation distance, for any measurable map $g : \mathbb{R}^d \times \mathcal{Z} \to \mathcal{Z}'$ taking values in any measurable space $\mathcal{Z}'$, we have

$$\|\mathcal{L}(g(L_n, \mathbf{z}_1)) - \mathcal{L}(g(L_n, \mathbf{z}_0))\|_{\mathrm{TV}} \le \|\mathcal{L}((L_n, \mathbf{z}_1)) - \mathcal{L}((L_n, \mathbf{z}_0))\|_{\mathrm{TV}}.$$

The inequality holds in particular for $g$ defined for all $(l, z)$ in $\mathbb{R}^d \times \mathcal{Z}$ by $g(l, z) = (F(l), z)$, from which we get

$$\Delta_n(P, \mathcal{A}) \le \|\mathcal{L}((L_n, \mathbf{z}_1)) - \mathcal{L}((L_n, \mathbf{z}_0))\|_{\mathrm{TV}} = \mathbb{E}\left[\|\mathcal{L}(L_n \mid \mathbf{z}_1) - \mathcal{L}(L_n)\|_{\mathrm{TV}}\right],$$

in which the expectation is taken over $\mathbf{z}_1$.

For $j = 1, \ldots, n$, denote by $\mathbf{v}_j := C^{-1/2}(L(\mathbf{z}_j) - \mathbb{E}[L(\mathbf{z}_j)])$ the centered and reduced version of $L(\mathbf{z}_j)$. The total variation distance being invariant by translation and rescaling, we shall write

$$\|\mathcal{L}(L_n \mid \mathbf{z}_1) - \mathcal{L}(L_n)\|_{\mathrm{TV}} = \|\mathcal{L}(L_n - \mathbb{E}[L(\mathbf{z}_1)]) - \mathcal{L}(L_n - \mathbb{E}[L(\mathbf{z}_1)] \mid \mathbf{z}_1)\|_{\mathrm{TV}}$$

$$= \left\| \mathcal{L}\left( \frac{1}{n} \sum_{j=1}^n (L(\mathbf{z}_j) - \mathbb{E}[L(\mathbf{z}_j)]) \right) - \mathcal{L}\left( \frac{1}{n} \sum_{j=1}^n (L(\mathbf{z}_j) - \mathbb{E}[L(\mathbf{z}_j)]) \Big| \mathbf{z}_1 \right) \right\|_{\mathrm{TV}}$$

$$= \left\| \mathcal{L}\left( \frac{C^{-1/2}}{\sqrt{n}} \sum_{j=1}^n (L(\mathbf{z}_j) - \mathbb{E}[L(\mathbf{z}_j)]) \right) - \mathcal{L}\left( \frac{C^{-1/2}}{\sqrt{n}} \sum_{j=1}^n (L(\mathbf{z}_j) - \mathbb{E}[L(\mathbf{z}_j)]) \Big| \mathbf{z}_1 \right) \right\|_{\mathrm{TV}}$$

$$= \left\| \mathcal{L}\left( \frac{1}{\sqrt{n}} \sum_{j=1}^n \mathbf{v}_j \right) - \mathcal{L}\left( \frac{1}{\sqrt{n}} \sum_{j=1}^n \mathbf{v}_j \Big| \mathbf{v}_1 \right) \right\|_{\mathrm{TV}}.$$

Denoting by $\mathcal{N}_d(\beta, \Sigma)$ the $d-$dimensional normal distribution with parameters $(\beta, \Sigma)$, it holds almost surely that

$$\left\| \mathcal{L}\left( \frac{1}{\sqrt{n}} \sum_{j=1}^n \mathbf{v}_j \right) - \mathcal{L}\left( \frac{1}{\sqrt{n}} \sum_{j=1}^n \mathbf{v}_j \Big| \mathbf{v}_1 \right) \right\|_{\mathrm{TV}} \le \left\| \mathcal{L}\left( \frac{1}{\sqrt{n}} \sum_{j=1}^n \mathbf{v}_j \right) - \mathcal{N}_d(0, \mathbf{I}_d) \right\|_{\mathrm{TV}}$$

$$+ \left\| \mathcal{L}\left( \frac{1}{\sqrt{n}} \sum_{j=1}^n \mathbf{v}_j \Big| \mathbf{v}_1 \right) - \mathcal{N}_d\left( \frac{1}{\sqrt{n}} \mathbf{v}_1, \frac{n-1}{n} I_d \right) \right\|_{\mathrm{TV}}$$

$$+ \left\| \mathcal{N}_d(0, \mathbf{I}_d) - \mathcal{N}_d\left( \frac{1}{\sqrt{n}} \mathbf{v}_1, \frac{n-1}{n} \mathbf{I}_d \right) \right\|_{\mathrm{TV}}$$

$$= \left\| \mathcal{L}\left( \frac{1}{\sqrt{n}} \sum_{j=1}^n \mathbf{v}_j \right) - \mathcal{N}_d(0, \mathbf{I}_d) \right\|_{\mathrm{TV}}$$

$$+ \left\| \mathcal{L}\left( \frac{1}{\sqrt{n-1}} \sum_{j=1}^{n-1} \mathbf{v}_j \right) - \mathcal{N}_d(0, \mathbf{I}_d) \right\|_{\mathrm{TV}}$$

$$+ \left\| \mathcal{N}_d(0, \mathbf{I}_d) - \mathcal{N}_d\left( \frac{1}{\sqrt{n}} \mathbf{v}_1, \frac{n-1}{n} \mathbf{I}_d \right) \right\|_{\mathrm{TV}}.$$

Applying Theorem 2.6 of Bally and Caramellino [2016] with variable $\mathbf{v}_j$ and parameter $r = 2$, one can upper bound the first two terms by some constant $C(d)(1 + m_3)$ times $n^{-1/2}$. The constant $C(d)$ here depends only on the dimension of the parameters $d$. We may upper bound the last term by the following proposition

**Proposition G.1.** *Let $n$ be an integer and $\beta \in \mathbb{R}^d$ be any $d-$dimensional vector. Then it holds that*

$$\left\| \mathcal{N}_d(0, \mathbf{I}_d) - \mathcal{N}_d\left( \frac{1}{\sqrt{n}} \beta, \frac{n-1}{n} \mathbf{I}_d \right) \right\|_{TV} \le \frac{\sqrt{d}}{2n} + \frac{1}{2\sqrt{n}} \|\beta\|_2.$$

Applying Proposition G.1 to the last quantity, it holds that

$$\left\| \mathcal{N}_d(0, \mathbf{I}_d) - \mathcal{N}_d\left( \frac{1}{\sqrt{n}} \mathbf{v}_1, \frac{n-1}{n} \mathbf{I}_d \right) \right\|_{\mathrm{TV}} \le \frac{\sqrt{d}}{2n} + \frac{1}{2\sqrt{n}} \|\mathbf{v}_1\|_2,$$

and the result follows from taking the expectation, with $c_{L,P} = C(d)(1 + m_3) + \frac{m_1}{2}$. $\qquad\square$

*Proof of Proposition G.1.* Applying Proposition 2.1 of Devroye et al. [2018], it holds almost surely that

$$\left\| \mathcal{N}_d(0, \boldsymbol{I}_d) - \mathcal{N}_d\left(\frac{1}{\sqrt{n}}\beta, \frac{n-1}{n}\boldsymbol{I}_d\right) \right\|_{\text{TV}}$$

$$\leq \frac{1}{2}\sqrt{tr\left(\boldsymbol{I}_d\frac{n-1}{n}\boldsymbol{I}_d - \boldsymbol{I}_d\right) + \frac{1}{n}\|\beta\|_2^2 - \ln\left(det\left(\frac{n-1}{n}\boldsymbol{I}_d\right)\right)}$$

$$= \frac{1}{2}\sqrt{-\frac{d}{n} + \frac{1}{n}\|\beta\|_2^2 - d\ln\left(\frac{n-1}{n}\right)}$$

$$\leq \frac{1}{2}\sqrt{-d\left(\frac{1}{n} + \ln\left(\frac{n-1}{n}\right)\right)} + \frac{1}{2}\sqrt{\frac{1}{n}\|\beta\|_2^2}$$

$$\leq \frac{\sqrt{d}}{2n} + \frac{1}{2\sqrt{n}}\|\beta\|_2,$$

where $tr(\cdot)$ is the trace operator and $det(\cdot)$ is the matrix determinant operator. The third inequality is due to $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$ for positive scalars $a$ and $b$. The first term in the last inequality comes from the fact that $x - 1 - \ln(x) \leq (x-1)^2$ if $x \geq 1/3$ which holds with $x = \frac{n-1}{n}$ for $n \geq 2$.

□

*Proof of Remark 5.1.* Setting $c := C(d)(1 + m_3) + \frac{m_1}{2}$, from Equation 9, we have that

$$cn^{-1/2} + \frac{\sqrt{d}}{2}n^{-1} \leq \varepsilon,$$

is sufficient to ensure $\Delta_n(P, \mathcal{A}) \leq \varepsilon$, hence a security of at least $1 - \varepsilon$.

Setting $x := n^{-1/2}$, it is equivalent to

$$cx + \frac{\sqrt{d}}{2}x^2 - \varepsilon \leq 0.$$

From the the study of the above quadratic function, as $x \geq 0$ is assumed, we get that this is equivalent to

$$n^{-1/2} \leq \frac{-c + \sqrt{c^2 + 2\varepsilon\sqrt{d}}}{\sqrt{d}}$$

$$\iff n \geq \frac{d}{2c^2 + 2\varepsilon\sqrt{d} - 2c\sqrt{c^2 + 2\varepsilon\sqrt{d}}}$$

$$= \frac{d}{2c^2} \frac{1}{1 + \frac{\varepsilon\sqrt{d}}{c^2} - \sqrt{1 + 2\frac{\varepsilon\sqrt{d}}{c^2}}}.$$

From the mean-value form of Taylor theorem of order 2 at 0, there exists $0 \leq \bar{u} \leq u := \frac{\varepsilon\sqrt{d}}{c^2}$ such that

$$\sqrt{1 + 2u} = 1 + u - \frac{1}{2}(1 + 2\bar{u})^{-3/2}.$$

Therefore, the condition becomes

$$n \geq \frac{d}{2c^2} \frac{2(1 + 2\bar{u})^{3/2}}{u^2}$$

$$= \varepsilon^{-2}c^2(1 + 2\bar{u})^{3/2}.$$

As $\bar{u} \leq u \leq \frac{\sqrt{d}}{c^2}$, $n \geq \varepsilon^{-2}c^2(1 + \frac{\sqrt{d}}{c^2})^{3/2}$ ensures the above condition, hence the result.

□

In all proofs below, we make use of the following fact

$$\Delta_n(P, \mathcal{A}) = \mathbb{E}\left[\left\|\mathcal{L}(\hat{\theta}_n) - \mathcal{L}(\hat{\theta}_n|z_1)\right\|_{\mathrm{TV}}\right], \tag{24}$$

where $\mathcal{L}(\hat{\theta}_n|z_1)$ is the distribution of $\hat{\theta}_n$ conditional to $z_1$, and the expectation is taken on the random variable $z_1$.

*Proof of Theorem 5.4.* The proof will be divided in two steps. First, we will prove the inequality

$$\Delta_n(P, \mathcal{A}) \le \frac{1}{2} \sum_{j=1}^{K} \mathbb{E}\left[\left|\frac{B_j}{n} - p_j\right|\right], \tag{25}$$

for any distribution $P$ and algorithm $\mathcal{A}$. Second, we prove that this upper bound is reached for algorithms that map any data set to a Dirac mass, summarized in the following lemma.

**Lemma G.2.** *For $j = 1, \ldots, K$, let $B_j$ be random variables having Binomial distribution with parameters $(n, p_j)$. Suppose that $\mathcal{A}(z_1, \cdots, z_n) = \delta_{F\left(\frac{1}{n}\sum_{j=1}^{n}\delta_{z_j}\right)}$ for any $n \in \mathbb{N}$ and $z_1, \ldots, z_n \in \mathcal{Z}$, for some measurable map $F: \mathcal{M} \to \Theta$ with infinite range $|\Theta| = \infty$, i.e. $\hat{\theta}_n \overset{\mathcal{L}}{=} F\left(\frac{1}{n}\sum_{j=1}^{n}\delta_{z_j}\right)$. Then we have*

$$\max_F \Delta_n(P, \mathcal{A}) = \frac{1}{2} \sum_{j=1}^{K} \mathbb{E}\left[\left|\frac{B_j}{n} - p_j\right|\right].$$

Theorem 5.4 will simply follow from Lemma G.2 and Equation 25.

Let us first prove Equation 25.
Since $\hat{\theta}_n$ has distribution $G(\hat{P}_n)$ conditionally on $\mathbf{z}$, where $\hat{P}_n := \frac{1}{n}\sum_{j=1}^{n}\delta_{z_j}$ is the empirical distribution of the data set, from Proposition A.3, we have

$$P(\hat{\theta}_n \in B) = \mathbb{E}[P(\hat{\theta}_n \in B|\mathbf{z})]$$
$$= \mathbb{E}[G(\hat{P}_n)(B)] \tag{26}$$
$$P(\hat{\theta}_n \in B|z_1) = \mathbb{E}[G(\hat{P}_n)(B)|z_1], \tag{27}$$

for any measurable set $B$.
Recall that $u_1, \ldots, u_K$ are the (fixed) support points of $P$. For any $k \in \{1, \cdots, K\}$, let $\hat{P}_n^k := \frac{1}{n}\left(\delta_{u_k} + \sum_{j=2}^{n}\delta_{z_j}\right)$. Using Equations 24, 26 and 27 we may rewrite $\Delta_n(P, \mathcal{A})$ as

$$\Delta_n(P, \mathcal{A}) = \sum_{k=1}^{K} p_k \sup_B \left(\mathbb{E}[G(\hat{P}_n)(B)] - \mathbb{E}[G(\hat{P}_n^k)(B)]\right). \tag{28}$$

For any integer $n$, let $\mathcal{M}_n$ be the set of all possible empirical distributions for data sets with $n$ points and let $\mathcal{G}_n = G(\mathcal{M}_n)$. Since $P$ has at most countable support, then $\mathcal{G}_n$ is at most countable and Equation 28 gives

$$\Delta_n(P, \mathcal{A}) = \sum_{k=1}^{K} p_k \sup_B \left(\sum_{g \in \mathcal{G}_n} g(B) P(G(\hat{P}_n) = g) - \sum_{g \in \mathcal{G}_n} g(B) P(G(\hat{P}_n^k) = g)\right). \tag{29}$$

For some fixed $g \in \mathcal{G}_n$, let us denote by $\mathcal{M}_n(g) = G^{-1}(\{g\}) \cap \mathcal{M}_n$ the set of possible empirical distributions $Q$ in $\mathcal{M}_n$ such that $G(Q) = g$. Then we have for any $g \in \mathcal{G}_n$,

$$g(B)\left(P(G(\hat{P}_n) = g) - P(G(\hat{P}_n^k) = g)\right) = \sum_{Q \in \mathcal{M}_n(g)} G(Q)(B)\left(P(\hat{P}_n = Q) - P(\hat{P}_n^k = Q)\right),$$

so that summing over all $g$ gives

$$\Delta_n(P, \mathcal{A}) = \sum_{k=1}^{K} p_k \sup_B \left( \sum_{Q \in \mathcal{M}_n} G(Q)(B) \left[ P(\hat{P}_n = Q) - P(\hat{P}_n^k = Q) \right] \right), \tag{30}$$

since $(\mathcal{M}_n(g))_{g \in \mathcal{G}_n}$ is a partition of $\mathcal{M}_n$. As the distribution is discrete, any possible value $Q$ of $\hat{P}_n$ is uniquely determined by a $K-$tuple $(k_1, \cdots, k_K)$ (if $K = \infty$ then by a sequence $(k_1, k_2, \cdots)$) of non-negative integers such that $\sum_{j=1}^{K} k_j = n$ and $Q = \frac{1}{n} \sum_{j=1}^{K} k_j \delta_{u_j}$. The $K-$tuple (or sequence) corresponds to the distribution of the samples among the atoms, that is, if we define, for $j = 1, \dots, K$, the random variable $N_j$ as the number of samples in the dataset equal to $u_j$, then for such $Q$,

$$P(\hat{P}_n = Q) = P(N_j = k_j; \ j = 1, \dots, K).$$

Since the samples are i.i.d., we get for such $Q$

$$P(\hat{P}_n = Q) = \binom{n}{k_1, \cdots, k_K} \prod_{j=1}^{K} p_j^{k_j}, \tag{31}$$

where $\binom{n}{k_1, \cdots, k_m} = \frac{n!}{k_1! \cdots k_m!}$ is the multinomial coefficient. Notice that when $K = +\infty$, only a finite number $m$ of integers $k_j$ are non zero, so that Equation 31 can be understood to hold also when $K = +\infty$ by keeping only the terms involving the positive integers $k_j$.

Let us now compute $P(\hat{P}_n^1 = Q)$. If $k_1 = 0$, then $P(\hat{P}_n^1 = Q) = 0$. Else,

$$P(\hat{P}_n^1 = Q) = P(N_1 = k_1 - 1, \ N_j = k_j; \ j = 2, \dots, K)$$

$$= \binom{n-1}{k_1 - 1, k_2, \cdots, k_K} \left( \prod_{j=2}^{K} p_j^{k_j} \right) p_1^{k_1 - 1}$$

$$= \frac{k_1}{np_1} \binom{n}{k_1, \cdots, k_K} \prod_{j=1}^{K} p_j^{k_j},$$

which again is understood to hold also when $K = +\infty$.
Therefore in both cases, we get

$$P(\hat{P}_n^1 = Q) = \frac{k_1}{np_1} \binom{n}{k_1, \cdots, k_K} \prod_{j=1}^{K} p_j^{k_j}. \tag{32}$$

Now, using 31 and 32, denoting by $g_N$ the image by $G$ of the distribution determined by the $K-$tuple $N = (k_1, \cdots, k_K)$, we get

$$\sum_{Q \in \mathcal{M}_n} G(Q)(B) \left( P(\hat{P}_n = Q) - P(\hat{P}_n^1 = Q) \right) = \sum_{k_1 + \cdots + k_K = n} g_N(B) \binom{n}{k_1, \cdots, k_K} \prod_{j=1}^{K} p_j^{k_j} \left( 1 - \frac{k_1}{np_1} \right)$$

$$= \mathbb{E} \left[ \left( 1 - \frac{N_1}{np_1} \right) g_N(B) \right],$$

where $N = (N_1, \cdots, N_K)$ follows a multinomial distribution of parameters $(n; p_1, \cdots, p_K)$. The computation being similar for any $k = 1, \dots, K$, we easily obtain that for any $k = 1, \dots, K$

$$\sum_{Q \in \mathcal{M}_n} G(Q)(B) \left( P(\hat{P}_n = Q) - P(\hat{P}_n^k = Q) \right) = \mathbb{E} \left[ \left( 1 - \frac{N_k}{np_k} \right) g_N(B) \right],$$

Now, plugging it into Equation 30 gives

$$\Delta_n(P, \mathcal{A}) = \sum_{k=1}^{K} p_k \sup_B \mathbb{E} \left[ \left( 1 - \frac{N_k}{np_k} \right) g_N(B) \right]. \tag{33}$$

For any real number $x \in \mathbb{R}$, we denote by $(x)_+ = \max(x, 0)$ its positive part and $(x)_- = \max(0, -x)$ its negative part. We get from Equation 33

$$\Delta_n(P, \mathcal{A}) \leq \sum_{k=1}^{K} p_k \mathbb{E}\left[\sup_B \left(1 - \frac{N_k}{np_k}\right) g_N(B)\right]$$
$$= \sum_{k=1} p_k \mathbb{E}\left[\left(1 - \frac{N_k}{np_k}\right)_+\right]$$
(34)

$$\Delta_n(P, \mathcal{A}) \leq \sum_{k=1} p_k \mathbb{E}\left[\left(1 - \frac{N_k}{np_k}\right)_-\right]$$
(35)

where the equality in Equation 34)comes from the fact that the supremum is reached on null sets when $1 - N_k/np_k$ is negative, and on sets of mass 1 when it is positive. Equation 35 is obtained by replacing $B$ by its complementary $B^c$ in the supremum and remarking that $\mathbb{E}[1 - N_k/np_k] = 0$. Combining Equations 34 and 35 gives

$$\Delta_n(P; \mathcal{A}) \leq \sum_{k=1}^{K} p_k \min\left\{\mathbb{E}\left[\left(1 - \frac{N_k}{np_k}\right)_+\right]; \mathbb{E}\left[\left(1 - \frac{N_k}{np_k}\right)_-\right]\right\}$$
$$\leq \frac{1}{2} \sum_{k=1}^{K} \mathbb{E}\left[\left|1 - \frac{N_k}{np_k}\right|\right],$$

which proves Equation 25.

$\square$

*Proof of Lemma G.2.* For some fixed $\theta \in \Theta$, we similarly denote by $\mathcal{M}_n(\theta) = F^{-1}(\{\theta\}) \cap \mathcal{M}_n$ the set of possible empirical distributions $Q$ in $\mathcal{M}_n$ such that $F(Q) = \theta$. Using Equation 24, and following similar steps as in Equations 28, 29 and 30, by triangular inequality, we get that

$$\Delta_n(P, \mathcal{A}) = \sum_{k=1}^{K} \frac{p_k}{2} \sum_{g \in \mathcal{G}_n} \left|P(\delta_{F(\hat{P}_n)} = g) - P(\delta_{F(\hat{P}_n^k)} = g)\right|$$
$$= \sum_{k=1}^{K} \frac{p_k}{2} \sum_{\theta \in \Theta} \left|P(F(\hat{P}_n) = \theta) - P(F(\hat{P}_n^k) = \theta)\right|$$
$$= \sum_{k=1}^{K} \frac{p_k}{2} \sum_{\theta \in \Theta} \left|\sum_{Q \in \mathcal{M}_n(\theta)} \left(P(\hat{P}_n = Q) - P(\hat{P}_n^k = Q)\right)\right|$$
(36)
$$\leq \sum_{k=1}^{K} \frac{p_k}{2} \sum_{Q \in \mathcal{M}_n} \left|P(\hat{P}_n = Q) - P(\hat{P}_n^k = Q)\right|,$$

since $(\mathcal{M}_n(\theta))_{\theta \in \Theta}$ is a partition of $\mathcal{M}_n$. We now prove that when taking the maximum over all possible measurable maps $F$ having range $\Theta$, the inequality becomes an equality. Indeed, since $\Theta$ is infinite, it is possible to construct $F$ such that $F$ is an injection from $\bigcup_{n \in \mathbb{N}} \mathcal{M}_n$ to $\Theta$, in which case for all $\theta \in \Theta$, $\mathcal{M}_n(\theta)$ is either the emptyset or a singleton. Thus, Equation 36 gives

$$\max_F \Delta_n(P, \mathcal{A}) = \sum_{k=1}^{K} \frac{p_k}{2} \sum_{Q \in \mathcal{M}_n} \left|P(\hat{P}_n = Q) - P(\hat{P}_n^k = Q)\right|,$$

and the lemma follows from Equations 31 and 32 and the same steps as in the proof of Theorem 5.4.

$\square$

*Proof of Corollary 5.4.1.* The upper bound comes from Cauchy-Schwartz's inequality and Theorem 5.4 since for any $j = 1, \ldots, K$,

$$\mathbb{E}\left[\left|\frac{B_j}{n} - p_j\right|\right] \leq \sqrt{Var(B_j/n)} = \sqrt{p_j(1-p_j)}n^{-1/2}.$$

We now prove the lower bound.

Define $m_k := \lfloor np_k \rfloor$, $k = 1, \ldots, K$. Using Theorem 5.4 , and De Moivre [1730], it holds that

$$\max_{\mathcal{A}} \Delta_n(P, \mathcal{A}) = \frac{1}{n} \sum_{k=1}^{K} \binom{n}{m_k + 1}(m_k + 1)p_k^{m_k+1}(1 - p_k)^{n-m_k}, \tag{37}$$

where $\binom{n}{m_k+1} = \frac{n!}{(m_k+1)!(n-(m_k+1))!}$ is a binomial coefficient. We shall approximate this binomial coefficient by Robbins [1955], which states that for any integer $k \geq 1$, it holds that

$$\sqrt{2\pi}k^{k+1/2}e^{-k}e^{1/12(k+1)} < k! < \sqrt{2\pi}k^{k+1/2}e^{-k}e^{1/12k}. \tag{38}$$

Note first that since for any $k = 1, \cdots, K$, $n > 1/p_k$, then also for any $k = 1, \cdots, K$, $n > 1/(1 - p_k)$. This implies that for any $k = 1, \cdots, K$,

$$1 \leq m_k \leq n - 2. \tag{39}$$

Set $a_k := \exp\left(\frac{1}{12(n+1)} - \frac{1}{12(m_k+1)} - \frac{1}{12(n-(m_k+1))}\right)$. Using equation 39 we get

$$a_k \geq \exp(-1/6). \tag{40}$$

One may apply Inequality 38 to get

$$\binom{n}{m_k + 1} \overset{38}{>} \frac{\sqrt{2\pi}n^{n+1/2}}{\sqrt{2\pi}(m_k+1)^{m_k+1+1/2}\sqrt{2\pi}(n-(m_k+1))^{n-(m_k+1)+1/2}} \frac{e^{-n}}{e^{-(m_k+1)}e^{-(n-(m_k+1))}} a_k$$

$$= \frac{\sqrt{n}}{\sqrt{2\pi}}n^n \left[m_k + 1\right]^{-(m_k+1+1/2)} \left[n - (m_k + 1)\right]^{-(n-(m_k+1)+1/2)} a_k$$

$$:= c_k a_k.$$

Now,

$$c_k(m_k + 1)p_k^{m_k+1}(1 - p_k)^{n-m_k} = \frac{\sqrt{n}}{\sqrt{2\pi}}n^n \left[m_k + 1\right]^{-(m_k+1+1/2-1)} \left[n - (m_k + 1)\right]^{-(n-(m_k+1)+1/2)}$$

$$\times p_k^{m_k+1}(1 - p_k)^{n-m_k}$$

$$= \frac{\sqrt{n}}{\sqrt{2\pi}}n^n (np_k)^{-(m_k+1/2)} \left[\frac{m_k + 1}{np_k}\right]^{-(m_k+1/2)}$$

$$\times (n(1 - p_k))^{-(n-(m_k+1/2))} \left[\frac{n - (m_k + 1)}{n(1 - p_k)}\right]^{-(n-(m_k+1/2))}$$

$$\times p_k^{m_k+1}(1 - p_k)^{n-m_k}$$

$$= \frac{\sqrt{n}}{\sqrt{2\pi}}\sqrt{p_k(1 - p_k)} \left[\frac{m_k + 1}{np_k}\right]^{-(m_k+1/2)}$$

$$\times \left[\frac{n - (m_k + 1)}{n(1 - p_k)}\right]^{-(n-(m_k+1/2))}$$

$$:= \frac{\sqrt{n}}{\sqrt{2\pi}}\sqrt{p_k(1 - p_k)}d_k,$$

which finally implies

$$\frac{\sqrt{n}}{\sqrt{2\pi}}\sqrt{p_k(1-p_k)}d_k a_k < \binom{n}{m_k+1}(m_k+1)p_k^{m_k+1}(1-p_k)^{n-m_k}.$$

Define $\epsilon_k \in [0,1)$ such that $np_k = m_k + \epsilon_k$. Then

$$d_k = \exp\left\{\left(m_k + \frac{1}{2}\right)\ln\left(\frac{m_k+\epsilon_k}{m_k+1}\right) + \left(n-m_k-\frac{1}{2}\right)\ln\left(\frac{n-m_k-\epsilon_k}{n-m_k-1}\right)\right\}$$

For any $m \in \{1; \ldots; n-2\}$, $\epsilon \in [0,1)$, define

$$f(m,\epsilon) = \left(m+\frac{1}{2}\right)\ln\left(\frac{m+\epsilon}{m+1}\right) + \left(n-m-\frac{1}{2}\right)\ln\left(\frac{n-m-\epsilon}{n-m-1}\right).$$

By studying the function $\epsilon \mapsto f(m,\epsilon)$ we get that for all $\epsilon \in [0,1)$, $f(m,\epsilon) \geq \min\{f(m,0),0\}$. By studying the function $m \mapsto f(m,0)$ we get that for all $m \in \{1; \ldots; n-2\}$, $f(m,0) \geq \min\{f(1,0), f(n-2,0)\}$. But

$$f(1,0) = -f(n-2,0) = -\frac{3}{2}\log(2) + \left(n-\frac{3}{2}\right)\log\left(1+\frac{1}{n-2}\right).$$

Now, Taylor expansion of $\log(1+u)$ allows to prove

$$-\frac{3}{2}\log(2) + 1 - \frac{1}{4(n-2)^2} \leq f(1,0) \leq -\frac{3}{2}\log(2) + 1 + \frac{1}{2(n-2)}$$

When $n \geq 5$, it is easy to see that $-\frac{3}{2}\log(2) + 1 - \frac{1}{4(n-2)^2} > 0$, so that using Equation 40, we get the result with

$$c = \frac{\exp\left(\frac{3}{2}\log(2) - 1 - 1/3\right)}{\sqrt{2\pi}}.$$

A rough approximation gives $c > 0.29$. $\qquad\square$