

## SOME THEORETICAL PROPERTIES OF GANS

BY GÉRARD BIAU<sup>\*</sup>, BENOÎT CADRE<sup>†</sup>, MAXIME SANGNIER<sup>\*</sup> AND UGO  
TANIELIAN<sup>\*,‡</sup>

*Sorbonne Université<sup>\*</sup>, Univ Rennes<sup>†</sup> and Criteo AI Lab<sup>‡</sup>*

Generative Adversarial Networks (GANs) are a class of generative algorithms that have been shown to produce state-of-the-art samples, especially in the domain of image creation. The fundamental principle of GANs is to approximate the unknown distribution of a given data set by optimizing an objective function through an adversarial game between a family of generators and a family of discriminators. In this paper, we offer a better theoretical understanding of GANs by analyzing some of their mathematical and statistical properties. We study the deep connection between the adversarial principle underlying GANs and the Jensen-Shannon divergence, together with some optimality characteristics of the problem. An analysis of the role of the discriminator family via approximation arguments is also provided. In addition, taking a statistical point of view, we study the large sample properties of the estimated distribution and prove in particular a central limit theorem. Some of our results are illustrated with simulated examples.

**1. Introduction.** The fields of machine learning and artificial intelligence have seen spectacular advances in recent years, one of the most promising being perhaps the success of Generative Adversarial Networks (GANs), introduced by Goodfellow et al. (2014). GANs are a class of generative algorithms implemented by a system of two neural networks contesting with each other in a zero-sum game framework. This technique is now recognized as being capable of generating photographs that look authentic to human observers (e.g., Salimans et al., 2016), and its spectrum of applications is growing at a fast pace, with impressive results in the domains of inpainting, speech, and 3D modeling, to name but a few. A survey of the most recent advances is given by Goodfellow (2016).

The objective of GANs is to generate fake observations of a target distribution  $p^*$  from which only a true sample (e.g., real-life images represented using raw pixels) is available. It should be pointed out at the outset that the data involved in the domain are usually so complex that no exhaustive description of  $p^*$  by a classical parametric model is appropriate, nor its estimation by a traditional maximum likelihood approach. Similarly, the dimension of the samples is often very large, and this effectively excludes a strategy based on nonparametric density estimation

---

*MSC 2010 subject classifications:* Primary 62F12; secondary 68T01

*Keywords and phrases:* generative models, adversarial principle, Jensen-Shannon divergence, neural networks, central limit theorem

techniques such as kernel or nearest neighbor smoothing, for example. In order to generate according to  $p^*$ , GANs proceed by an adversarial scheme involving two components: a family of generators and a family of discriminators, which are both implemented by neural networks. The generators admit low-dimensional random observations with a known distribution (typically Gaussian or uniform) as input, and attempt to transform them into fake data that can match the distribution  $p^*$ ; on the other hand, the discriminators aim to accurately discriminate between the true observations from  $p^*$  and those produced by the generators. The generators and the discriminators are calibrated by optimizing an objective function in such a way that the distribution of the generated sample is as indistinguishable as possible from that of the original data. In pictorial terms, this process is often compared to a game of cops and robbers, in which a team of counterfeiters illegally produces banknotes and tries to make them undetectable in the eyes of a team of police officers, whose objective is of course the opposite. The competition pushes both teams to improve their methods until counterfeit money becomes indistinguishable (or not) from genuine currency.

From a mathematical point of view, here is how the generative process of GANs can be represented. All the densities that we consider in the article are supposed to be dominated by a fixed, known, measure  $\mu$  on  $E$ , where  $E$  is a Borel subset of  $\mathbb{R}^d$ . Depending on the practical context, this dominating measure may be the Lebesgue measure, the counting measure, or more generally the Hausdorff measure on some submanifold of  $\mathbb{R}^d$ . We assume to have at hand an i.i.d. sample  $X_1, \dots, X_n$ , drawn according to some unknown density  $p^*$  on  $E$ . These random variables model the available data, such as images or video sequences; they typically take their values in a high-dimensional space, so that the ambient dimension  $d$  must be thought of as large. The generators as a whole have the form of a parametric family of functions from  $\mathbb{R}^{d'}$  to  $E$  (usually,  $d' \ll d$ ), say  $\mathcal{G} = \{G_\theta\}_{\theta \in \Theta}$ ,  $\Theta \subset \mathbb{R}^p$ . Each function  $G_\theta$  is intended to be applied to a  $d'$ -dimensional random variable  $Z$  (sometimes called the noise)—in most cases Gaussian or uniform), so that there is a natural family of densities associated with the generators, say  $\mathcal{P} = \{p_\theta\}_{\theta \in \Theta}$ , where, by definition,  $G_\theta(Z) \stackrel{\mathcal{L}}{=} p_\theta d\mu$ . In this model, each density  $p_\theta$  is a potential candidate to represent  $p^*$ . On the other hand, the discriminators are described by a family of Borel functions from  $E$  to  $[0, 1]$ , say  $\mathcal{D}$ , where each  $D \in \mathcal{D}$  must be thought of as the probability that an observation comes from  $p^*$  (the higher  $D(x)$ , the higher the probability that  $x$  is drawn from  $p^*$ ). At some point, but not always, we will assume that  $\mathcal{D}$  is in fact a parametric class, of the form  $\{D_\alpha\}_{\alpha \in \Lambda}$ ,  $\Lambda \subset \mathbb{R}^q$ , as is always the case in practice. In GANs algorithms, both parametric models  $\{G_\theta\}_{\theta \in \Theta}$  and  $\{D_\alpha\}_{\alpha \in \Lambda}$  take the form of neural networks, but this does not play a fundamental role in this paper. We will simply remember that the dimensions  $p$  and  $q$  are potentially very large, which takes us away from a classical parametric setting. We also

insist on the fact that it is not assumed that  $p^*$  belongs to  $\mathcal{D}$ .

Let  $Z_1, \dots, Z_n$  be an i.i.d. sample of random variables, all distributed as the noise  $Z$ . The objective is to solve in  $\theta$  the problem

$$(1.1) \quad \inf_{\theta \in \Theta} \sup_{D \in \mathcal{D}} \left[ \prod_{i=1}^n D(X_i) \times \prod_{i=1}^n (1 - D \circ G_\theta(Z_i)) \right],$$

or, equivalently, to find  $\hat{\theta} \in \Theta$  such that

$$(1.2) \quad \sup_{D \in \mathcal{D}} \hat{L}(\hat{\theta}, D) \leq \sup_{D \in \mathcal{D}} \hat{L}(\theta, D), \quad \forall \theta \in \Theta,$$

where

$$\hat{L}(\theta, D) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n \ln D(X_i) + \frac{1}{n} \sum_{i=1}^n \ln(1 - D \circ G_\theta(Z_i))$$

( $\ln$  is the natural logarithm). The zero-sum game (1.1) is the statistical translation of making the distribution of  $G_\theta(Z_i)$  (i.e.,  $p_\theta$ ) as indistinguishable as possible from that of  $X_i$  (i.e.,  $p^*$ ). Here, distinguishability is understood as the capability to determine from which distribution an observation  $x$  comes from. Mathematically, this is captured by the discrimination value  $D(x)$ , which represents the probability that  $x$  comes from  $p^*$  rather than from  $p_\theta$ . Therefore, for a given  $\theta$ , the discriminator  $D$  is determined so as to be maximal on the  $X_i$  and minimal on the  $G_\theta(Z_i)$ . In the most favorable situation (that is, when the two samples are scattered by  $\mathcal{D}$ ),  $\sup_{D \in \mathcal{D}} \hat{L}(\theta, D)$  is zero, and the larger this quantity, the more distinguishable the two samples are. Hence, in order to make the distribution  $p_\theta$  as indistinguishable as possible from  $p^*$ ,  $G_\theta$  has to be driven so as to minimize  $\sup_{D \in \mathcal{D}} \hat{L}(\theta, D)$ .

As mentioned above, this adversarial problem is often illustrated by the struggle between a police team (the discriminators), trying to distinguish true banknotes from false ones (respectively, the  $X_i$  and the  $G_\theta(Z_i)$ ), and a counterfeiters team, slaving to produce banknotes as credible as possible and to mislead the police. Obviously, their objectives (represented by the quantity  $\hat{L}(\theta, D)$ ) are exactly opposite. All in all, we see that the criterion seeks to find the right balance between the conflicting interests of the generators and the discriminators. The hope is that the  $\hat{\theta}$  achieving equilibrium will make it possible to generate observations  $G_{\hat{\theta}}(Z_1), \dots, G_{\hat{\theta}}(Z_n)$  indistinguishable from reality, i.e., observations with a distribution close to the unknown  $p^*$ .

The criterion  $\hat{L}(\theta, D)$  involved in (1.2) is the criterion originally proposed in the adversarial framework of [Goodfellow et al. \(2014\)](#). Since then, the success of GANs in applications has led to a large volume of literature on variants, which all have many desirable properties but are based on different optimization criteria: examples are MMD-GANs ([Dziugaite et al., 2015](#)), f-GANs ([Nowozin et al., 2016](#)), Wasserstein-GANs ([Arjovsky et al., 2017](#)), and an approach based on scattering

transforms (Angles and Mallat, 2018). All these variations and their innumerable algorithmic versions constitute the galaxy of GANs. That being said, despite increasingly spectacular applications, little is known about the mathematical and statistical forces behind these algorithms (e.g., Arjovsky and Bottou, 2017; Liu et al., 2017; Liang, 2018; Zhang et al., 2018), and, in fact, nearly nothing about the primary adversarial problem (1.2). As acknowledged by Liu et al. (2017), basic questions on how well GANs can approximate the target distribution  $p^*$  remain largely unanswered. In particular, the role and impact of the discriminators on the quality of the approximation are still a mystery, and simple but fundamental questions regarding statistical consistency and rates of convergence remain open.

In the present article, we propose to take a small step towards a better theoretical understanding of GANs by analyzing some of the mathematical and statistical properties of the original adversarial problem (1.2). In Section 2, we study the deep connection between the population version of (1.2) and the Jensen-Shannon divergence, together with some optimality characteristics of the problem, often referred to in the literature but in fact poorly understood. Section 3 is devoted to a better comprehension of the role of the discriminator family via approximation arguments. Finally, taking a statistical point of view, we study in Section 4 the large sample properties of the distribution  $p_{\hat{\theta}}$  and of  $\hat{\theta}$ , and prove in particular a central limit theorem for this parameter. Section 5 summarizes the main results and discusses research directions for future work. For clarity, most technical proofs are gathered in Section 6. Some of our results are illustrated with simulated examples.

**2. Optimality properties.** We start by studying some important properties of the adversarial principle, emphasizing the role played by the Jensen-Shannon divergence. We recall that if  $P$  and  $Q$  are probability measures on  $E$ , and  $P$  is absolutely continuous with respect to  $Q$ , then the Kullback-Leibler divergence from  $Q$  to  $P$  is defined as

$$D_{\text{KL}}(P \parallel Q) = \int \ln \frac{dP}{dQ} dP,$$

where  $\frac{dP}{dQ}$  is the Radon-Nikodym derivative of  $P$  with respect to  $Q$ . The Kullback-Leibler divergence is always nonnegative, with  $D_{\text{KL}}(P \parallel Q)$  zero if and only if  $P = Q$ . If  $p = \frac{dP}{d\mu}$  and  $q = \frac{dQ}{d\mu}$  exist (meaning that  $P$  and  $Q$  are absolutely continuous with respect to  $\mu$ , with densities  $p$  and  $q$ ), then the Kullback-Leibler divergence is given as

$$D_{\text{KL}}(P \parallel Q) = \int p \ln \frac{p}{q} d\mu,$$

and alternatively denoted by  $D_{\text{KL}}(p \parallel q)$ . We also recall that the Jensen-Shannon divergence is a symmetrized version of the Kullback-Leibler divergence. It is de-

finned for any probability measures  $P$  and  $Q$  on  $E$  by

$$D_{\text{JS}}(P, Q) = \frac{1}{2}D_{\text{KL}}\left(P \parallel \frac{P+Q}{2}\right) + \frac{1}{2}D_{\text{KL}}\left(Q \parallel \frac{P+Q}{2}\right),$$

and satisfies  $0 \leq D_{\text{JS}}(P, Q) \leq \ln 2$ . The square root of the Jensen-Shannon divergence is a metric often referred to as Jensen-Shannon distance (Endres and Schindelin, 2003). When  $P$  and  $Q$  have densities  $p$  and  $q$  with respect to  $\mu$ , we use the notation  $D_{\text{JS}}(p, q)$  in place of  $D_{\text{JS}}(P, Q)$ .

For a generator  $G_\theta$  and an arbitrary discriminator  $D \in \mathcal{D}$ , the criterion  $\hat{L}(\theta, D)$  to be optimized in (1.2) is but the empirical version of the probabilistic criterion

$$L(\theta, D) \stackrel{\text{def}}{=} \int \ln(D)p^* d\mu + \int \ln(1-D)p_\theta d\mu.$$

We assume for the moment that the discriminator class  $\mathcal{D}$  is not restricted and equals  $\mathcal{D}_\infty$ , the set of all Borel functions from  $E$  to  $[0, 1]$ . We note however that, for all  $\theta \in \Theta$ ,

$$0 \geq \sup_{D \in \mathcal{D}_\infty} L(\theta, D) \geq -\ln 2 \left( \int p^* d\mu + \int p_\theta d\mu \right) = -\ln 4,$$

so that  $\inf_{\theta \in \Theta} \sup_{D \in \mathcal{D}_\infty} L(\theta, D) \in [-\ln 4, 0]$ . Thus,

$$\inf_{\theta \in \Theta} \sup_{D \in \mathcal{D}_\infty} L(\theta, D) = \inf_{\theta \in \Theta} \sup_{D \in \mathcal{D}_\infty: L(\theta, D) > -\infty} L(\theta, D).$$

This identity points out the importance of discriminators such that  $L(\theta, D) > -\infty$ , which we call  $\theta$ -admissible. In the sequel, in order to avoid unnecessary problems of integrability, we only consider such discriminators, keeping in mind that the others have no interest.

Of course, working with  $\mathcal{D}_\infty$  is somehow an idealized vision, since in practice the discriminators are always parameterized by some parameter  $\alpha \in \Lambda$ ,  $\Lambda \subset \mathbb{R}^q$ . Nevertheless, this point of view is informative and, in fact, is at the core of the connection between our generative problem and the Jensen-Shannon divergence. Indeed, taking the supremum of  $L(\theta, D)$  over  $\mathcal{D}_\infty$ , we have

$$\begin{aligned} \sup_{D \in \mathcal{D}_\infty} L(\theta, D) &= \sup_{D \in \mathcal{D}_\infty} \int [\ln(D)p^* + \ln(1-D)p_\theta] d\mu \\ &\leq \int \sup_{D \in \mathcal{D}_\infty} [\ln(D)p^* + \ln(1-D)p_\theta] d\mu \\ &= L(\theta, D_\theta^*), \end{aligned}$$

where

$$(2.1) \quad D_\theta^* \stackrel{\text{def}}{=} \frac{p^*}{p^* + p_\theta}.$$

(We use throughout the convention  $0/0 = 0$  and  $\infty \times 0 = 0$ .) By observing that  $L(\theta, D_\theta^*) = 2D_{\text{JS}}(p^*, p_\theta) - \ln 4$ , we conclude that, for all  $\theta \in \Theta$ ,

$$\sup_{D \in \mathcal{D}_\infty} L(\theta, D) = L(\theta, D_\theta^*) = 2D_{\text{JS}}(p^*, p_\theta) - \ln 4.$$

We note in particular that  $D_\theta^*$  is  $\theta$ -admissible. The fact that  $D_\theta^*$  realizes the supremum of  $L(\theta, D)$  over  $\mathcal{D}_\infty$  and that this supremum is connected to the Jensen-Shannon divergence between  $p^*$  and  $p_\theta$  appears in the original article by [Goodfellow et al. \(2014\)](#). This remark has given rise to many developments that interpret the adversarial problem (1.2) as the empirical version of the minimization problem  $\inf_\theta D_{\text{JS}}(p^*, p_\theta)$  over  $\Theta$ . Accordingly, many GANs algorithms try to learn the optimal function  $D_\theta^*$ , using for example stochastic gradient descent techniques and mini-batch approaches. However, it remains to prove that  $D_\theta^*$  is unique as a maximizer of  $L(\theta, D)$  over all  $D$ . The following theorem, which completes a result of ([Goodfellow et al., 2014](#)), shows that this is the case in some situations.

**THEOREM 2.1.** *Let  $\theta \in \Theta$  and  $D \in \mathcal{D}_\infty$  be such that  $L(\theta, D) = L(\theta, D_\theta^*)$ . Then  $D = D_\theta^*$   $\mu$ -almost everywhere on the complementary of the set  $\{p^* = p_\theta = 0\}$ . In particular, if  $\mu(\{p^* = p_\theta = 0\}) = 0$ , then the function  $D_\theta^*$  is the unique discriminator that achieves the supremum of the functional  $D \mapsto L(\theta, D)$  over  $\mathcal{D}_\infty$ , i.e.,*

$$\{D_\theta^*\} = \arg \max_{D \in \mathcal{D}_\infty} L(\theta, D).$$

Before proving the theorem, it is important to note that if we do not assume that  $\mu(\{p^* = p_\theta = 0\}) = 0$ , then we cannot conclude that  $D = D_\theta^*$   $\mu$ -almost everywhere. To see this, suppose that  $p_\theta = p^*$ . Then, whatever  $\bar{D} \in \mathcal{D}_\infty$  is, the discriminator  $D_\theta^* \mathbf{1}_{\{p_\theta > 0\}} + \bar{D} \mathbf{1}_{\{p_\theta = 0\}}$  satisfies

$$L(\theta, D_\theta^* \mathbf{1}_{\{p_\theta > 0\}} + \bar{D} \mathbf{1}_{\{p_\theta = 0\}}) = L(\theta, D_\theta^*).$$

This simple counterexample shows that uniqueness of the optimal discriminator does not hold in general.

**PROOF.** Let  $D \in \mathcal{D}_\infty$  be a discriminator such that  $L(\theta, D) = L(\theta, D_\theta^*)$ . In particular,  $L(\theta, D) > -\infty$  and  $D$  is  $\theta$ -admissible. Thus, letting  $A \stackrel{\text{def}}{=} \{p^* = p_\theta = 0\}$  and  $f_\alpha \stackrel{\text{def}}{=} p^* \ln(\alpha) + p_\theta \ln(1 - \alpha)$  for  $\alpha \in [0, 1]$ , we see that

$$\int_{A^c} (f_D - f_{D_\theta^*}) d\mu = 0.$$

Since, on  $A^c$ ,

$$f_D \leq \sup_{\alpha \in [0, 1]} f_\alpha = f_{D_\theta^*},$$

we have  $f_D = f_{D_\theta^*}$   $\mu$ -almost everywhere on  $A^c$ . By uniqueness of the maximizer of  $\alpha \mapsto f_\alpha$  on  $A^c$ , we conclude that  $D = D_\theta^*$   $\mu$ -almost everywhere on  $A^c$ .  $\square$

By definition of the optimal discriminator  $D_\theta^*$ , we have

$$L(\theta, D_\theta^*) = \sup_{D \in \mathcal{D}_\infty} L(\theta, D) = 2D_{\text{JS}}(p^*, p_\theta) - \ln 4, \quad \forall \theta \in \Theta.$$

Therefore, it makes sense to let the parameter  $\theta^* \in \Theta$  be defined as

$$L(\theta^*, D_{\theta^*}^*) \leq L(\theta, D_\theta^*), \quad \forall \theta \in \Theta,$$

or, equivalently,

$$(2.2) \quad D_{\text{JS}}(p^*, p_{\theta^*}) \leq D_{\text{JS}}(p^*, p_\theta), \quad \forall \theta \in \Theta.$$

The parameter  $\theta^*$  may be interpreted as the best parameter in  $\Theta$  for approaching the unknown density  $p^*$  in terms of Jensen-Shannon divergence, in a context where all possible discriminators are available. In other words, the generator  $G_{\theta^*}$  is the ideal generator, and the density  $p_{\theta^*}$  is the one we would ideally like to use to generate fake samples. Of course, whenever  $p^* \in \mathcal{P}$  (i.e., the target density is in the model), then  $p^* = p_{\theta^*}$ ,  $D_{\text{JS}}(p^*, p_{\theta^*}) = 0$ , and  $D_{\theta^*}^* = 1/2$ . This is, however, a very special case, which is of no interest, since in the applications covered by GANs, the data are usually so complex that the hypothesis  $p^* \in \mathcal{P}$  does not hold.

In the general case, our next theorem provides sufficient conditions for the existence and uniqueness of  $\theta^*$ . For  $P$  and  $Q$  probability measures on  $E$ , we let  $\delta(P, Q) = \sqrt{D_{\text{JS}}(P, Q)}$ , and recall that  $\delta$  is a distance on the set of probability measures on  $E$  (Endres and Schindelin, 2003). We let  $dP^* = p^*d\mu$  and, for all  $\theta \in \Theta$ ,  $dP_\theta = p_\theta d\mu$ .

**THEOREM 2.2.** *Assume that the model  $\mathcal{P} = \{P_\theta\}_{\theta \in \Theta}$  is convex and compact for the metric  $\delta$ . If  $p^* > 0$   $\mu$ -almost everywhere, then there exists a unique  $\bar{p} \in \mathcal{P}$  such that*

$$\{\bar{p}\} = \arg \min_{p \in \mathcal{P}} D_{\text{JS}}(p^*, p).$$

*In particular, if the model  $\mathcal{P}$  is identifiable, then*

$$\{\theta^*\} = \arg \min_{\theta \in \Theta} L(\theta, D_\theta^*)$$

*or, equivalently,*

$$\{\theta^*\} = \arg \min_{\theta \in \Theta} D_{\text{JS}}(p^*, p_\theta).$$

We note that the identifiability assumption in the second statement of the theorem is hardly satisfied in the high-dimensional context of (deep) neural networks. In this case, it is likely that several parameters  $\theta$  yield the same function (generator), so that the parametric setting is potentially misspecified. However, if we think in terms of distributions instead of parameters, then the first part of Theorem 2.2 ensures existence and uniqueness of the optimum.

PROOF. Assuming the first part of the theorem, the second one is obvious since  $L(\theta, D_\theta^*) = \sup_{D \in \mathcal{D}_\infty} L(\theta, D) = 2D_{\text{JS}}(p^*, p_\theta) - \ln 4$ . Therefore, it is enough to prove that there exists a unique density  $\bar{p}$  of  $\mathcal{P}$  such that

$$\{\bar{p}\} = \arg \min_{p \in \mathcal{P}} D_{\text{JS}}(p^*, p).$$

**Existence.** Since  $\mathcal{P}$  is compact for  $\delta$ , it is enough to show that the function

$$\begin{aligned} \mathcal{P} &\rightarrow \mathbb{R}_+ \\ P &\mapsto D_{\text{JS}}(P^*, P) \end{aligned}$$

is continuous. But this is clear since, for all  $P_1, P_2 \in \mathcal{P}$ ,  $|\delta(P^*, P_1) - \delta(P^*, P_2)| \leq \delta(P_1, P_2)$  by the triangle inequality. Therefore,  $\arg \min_{p \in \mathcal{P}} D_{\text{JS}}(p^*, p) \neq \emptyset$ .

**Uniqueness.** For  $a \geq 0$ , we consider the function  $F_a$  defined by

$$F_a(x) = a \ln \left( \frac{2a}{a+x} \right) + x \ln \left( \frac{2x}{a+x} \right), \quad x \geq 0,$$

with the convention  $0 \ln 0 = 0$ . Clearly,  $F_a''(x) = \frac{a}{x(a+x)}$ , which shows that  $F_a$  is strictly convex whenever  $a > 0$ . We now proceed to prove that  $L^1(\mu) \supset \mathcal{P} \ni p \mapsto D_{\text{JS}}(p^*, p)$  is strictly convex as well. Let  $\lambda \in (0, 1)$  and  $p_1, p_2 \in \mathcal{P}$  with  $p_1 \neq p_2$ , i.e.,  $\mu(\{p_1 \neq p_2\}) > 0$ . Then

$$\begin{aligned} &D_{\text{JS}}(p^*, \lambda p_1 + (1-\lambda)p_2) \\ &= \int F_{p^*}(\lambda p_1 + (1-\lambda)p_2) d\mu \\ &= \int_{\{p_1=p_2\}} F_{p^*}(p_1) d\mu + \int_{\{p_1 \neq p_2\}} F_{p^*}(\lambda p_1 + (1-\lambda)p_2) d\mu. \end{aligned}$$

By the strict convexity of  $F_{p^*}$  over  $\{p^* > 0\}$ , we obtain

$$\begin{aligned} &D_{\text{JS}}(p^*, \lambda p_1 + (1-\lambda)p_2) \\ &< \int_{\{p_1=p_2\}} F_{p^*}(p_1) d\mu + \lambda \int_{\{p_1 \neq p_2\}} F_{p^*}(p_1) d\mu + (1-\lambda) \int_{\{p_1 \neq p_2\}} F_{p^*}(p_2) d\mu, \end{aligned}$$



which implies

$$D_{\text{JS}}(p^*, \lambda p_1 + (1 - \lambda)p_2) < \lambda D_{\text{JS}}(p^*, p_1) + (1 - \lambda)D_{\text{JS}}(p^*, p_2).$$

Consequently, the function  $L^1(\mu) \supset \mathcal{P} \ni p \mapsto D_{\text{JS}}(p^*, p)$  is strictly convex, and its arg min over the convex set  $\mathcal{P}$  is either the empty set or a singleton.  $\square$

REMARK 2.1. *There are simple conditions for the model  $\mathcal{P} = \{P_\theta\}_{\theta \in \Theta}$  to be compact for the metric  $\delta$ . It is for example enough to suppose that  $\Theta$  is compact,  $\mathcal{P}$  is convex, and*

- (i) *For all  $x \in E$ , the function  $\theta \mapsto p_\theta(x)$  is continuous on  $\Theta$ ;*
- (ii) *One has  $\sup_{(\theta, \theta') \in \Theta^2} |p_\theta \ln p_{\theta'}| \in L^1(\mu)$ .*

*Let us quickly check that under these conditions,  $\mathcal{P}$  is compact for the metric  $\delta$ . Since  $\Theta$  is compact, by the sequential characterization of compact sets, it is enough to prove that if  $\Theta \supset (\theta_n)_n$  converges to  $\theta \in \Theta$ , then  $D_{\text{JS}}(p_\theta, p_{\theta_n}) \rightarrow 0$ . But,*

$$D_{\text{JS}}(p_\theta, p_{\theta_n}) = \int \left[ p_\theta \ln \left( \frac{2p_\theta}{p_\theta + p_{\theta_n}} \right) + p_{\theta_n} \ln \left( \frac{2p_{\theta_n}}{p_\theta + p_{\theta_n}} \right) \right] d\mu.$$

*By the convexity of  $\mathcal{P}$ , using (i) and (ii), the Lebesgue dominated convergence theorem shows that  $D_{\text{JS}}(p_\theta, p_{\theta_n}) \rightarrow 0$ , whence the result.*

Interpreting the adversarial problem in connection with the optimization program  $\inf_{\theta \in \Theta} D_{\text{JS}}(p^*, p_\theta)$  is a bit misleading, because this is based on the assumption that all possible discriminators are available (and in particular the optimal discriminator  $D_\theta^*$ ). In the end this means assuming that we know the distribution  $p^*$ , which is eventually not acceptable from a statistical perspective. In practice, the class of discriminators is always restricted to be a parametric family  $\mathcal{D} = \{D_\alpha\}_{\alpha \in \Lambda}$ ,  $\Lambda \subset \mathbb{R}^q$ , and it is with this class that we have to work. From our point of view, problem (1.2) is a likelihood-type problem involving two parametric families  $\mathcal{G}$  and  $\mathcal{D}$ , which must be analyzed as such, just as we would do for a classical maximum likelihood approach. In fact, it takes no more than a moment's thought to realize that the key lies in the approximation capabilities of the discriminator class  $\mathcal{D}$  with respect to the functions  $D_\theta^*$ ,  $\theta \in \Theta$ . This is the issue that we discuss in the next section.

**3. Approximation properties.** In the remainder of the article, we assume that  $\theta^*$  exists, keeping in mind that Theorem 2.2 provides us with precise conditions guaranteeing its existence and its uniqueness. As pointed out earlier, in practice only a parametric class  $\mathcal{D} = \{D_\alpha\}_{\alpha \in \Lambda}$ ,  $\Lambda \subset \mathbb{R}^q$ , is available, and it is therefore logical to consider the parameter  $\bar{\theta} \in \Theta$  defined by

$$\sup_{D \in \mathcal{D}} L(\bar{\theta}, D) \leq \sup_{D \in \mathcal{D}} L(\theta, D), \quad \forall \theta \in \Theta.$$

(We assume for now that  $\bar{\theta}$  exists—sufficient conditions for this existence, relating to compactness of  $\Theta$  and regularity of the model  $\mathcal{P}$ , will be given in the next section.) The density  $p_{\bar{\theta}}$  is thus the best candidate to imitate  $p_{\theta^*}$ , given the parametric families of generators  $\mathcal{G}$  and discriminators  $\mathcal{D}$ . The natural question is then: is it possible to quantify the proximity between  $p_{\bar{\theta}}$  and the ideal  $p_{\theta^*}$  via the approximation properties of the class  $\mathcal{D}$ ? In other words, if  $\mathcal{D}$  is growing, is it true that  $p_{\bar{\theta}}$  approaches  $p_{\theta^*}$ , and in the affirmative, in which sense and at which speed? Theorem 3.1 below provides a first answer to this important question, in terms of excess of Jensen-Shannon error  $D_{\text{JS}}(p^*, p_{\bar{\theta}}) - D_{\text{JS}}(p^*, p_{\theta^*})$ . To state the result, we will need an assumption.

Let  $\|\cdot\|_2$  be the  $L^2(\mu)$  norm. Our condition guarantees that the parametric class  $\mathcal{D}$  is rich enough to approach the discriminator  $D_{\bar{\theta}}^*$  in the  $L^2$  sense. In the remainder of the section, it is assumed that  $D_{\bar{\theta}}^* \in L^2(\mu)$ .

**Assumption ( $H_\varepsilon$ )** There exist  $\varepsilon > 0$ ,  $m \in (0, 1/2)$ , and  $D \in \mathcal{D} \cap L^2(\mu)$  such that  $m \leq D \leq 1 - m$  and  $\|D - D_{\bar{\theta}}^*\|_2 \leq \varepsilon$ .

We observe in passing that such a discriminator  $D$  is  $\bar{\theta}$ -admissible. We are now equipped to state our approximation theorem. For ease of reading, its proof is postponed to Section 6.

**THEOREM 3.1.** *Assume that, for some  $M > 0$ ,  $p^* \leq M$  and  $p_{\bar{\theta}} \leq M$ . Then, under Assumption ( $H_\varepsilon$ ) with  $\varepsilon < 1/(2M)$ , there exists a positive constant  $c$  (depending only upon  $m$  and  $M$ ) such that*

$$(3.1) \quad 0 \leq D_{\text{JS}}(p^*, p_{\bar{\theta}}) - D_{\text{JS}}(p^*, p_{\theta^*}) \leq c\varepsilon^2.$$

This theorem points out that if the class  $\mathcal{D}$  is rich enough to approximate the discriminator  $D_{\bar{\theta}}^*$  in such a way that  $\|D - D_{\bar{\theta}}^*\|_2 \leq \varepsilon$  for some small  $\varepsilon$ , then working with a restricted class of discriminators  $\mathcal{D}$  instead of the set of all discriminators  $\mathcal{D}_\infty$  has an impact that is not larger than a  $O(\varepsilon^2)$  factor with respect to the excess of Jensen-Shannon error. It shows in particular that the Jensen-Shannon divergence is a suitable criterion for the problem we are examining.

**4. Statistical analysis.** The data-dependent parameter  $\hat{\theta}$ , achieves the infimum of the adversarial problem (1.2). Practically speaking, it is this parameter that will be used in the end for producing fake data, via the associated generator  $G_{\hat{\theta}}$ . We first study in Subsection 4.1 the large sample properties of the distribution  $p_{\hat{\theta}}$  via the excess of Jensen-Shannon error  $D_{\text{JS}}(p^*, p_{\hat{\theta}}) - D_{\text{JS}}(p^*, p_{\theta^*})$ , and then state in Subsection 4.2 the almost sure convergence and asymptotic normality of the parameter  $\hat{\theta}$  as the sample size  $n$  tends to infinity. Throughout, the parameter sets  $\Theta$  and  $\Lambda$  are assumed to be compact subsets of  $\mathbb{R}^p$  and  $\mathbb{R}^q$ , respectively. To simplify

the analysis, we also assume that  $\mu(E) < \infty$ . In this case, every discriminator is in  $L^p(\mu)$  for all  $p \geq 1$ .

4.1. *Asymptotic properties of  $D_{\text{JS}}(p^*, p_{\hat{\theta}})$ .* As for now, we assume that we have at hand a parametric family of generators  $\mathcal{G} = \{G_\theta\}_{\theta \in \Theta}$ ,  $\Theta \subset \mathbb{R}^p$ , and a parametric family of discriminators  $\mathcal{D} = \{D_\alpha\}_{\alpha \in \Lambda}$ ,  $\Lambda \subset \mathbb{R}^q$ . We recall that the collection of probability densities associated with  $\mathcal{G}$  is  $\mathcal{P} = \{p_\theta\}_{\theta \in \Theta}$ , where  $G_\theta(Z) \stackrel{\mathcal{L}}{=} p_\theta d\mu$  and  $Z$  is some low-dimensional noise random variable. In order to avoid any confusion, for a given discriminator  $D = D_\alpha$  we use the notation  $\hat{L}(\theta, \alpha)$  (respectively,  $L(\theta, \alpha)$ ) instead of  $\hat{L}(\theta, D)$  (respectively,  $L(\theta, D)$ ) when useful. So,

$$\hat{L}(\theta, \alpha) = \frac{1}{n} \sum_{i=1}^n \ln D_\alpha(X_i) + \frac{1}{n} \sum_{i=1}^n \ln(1 - D_\alpha \circ G_\theta(Z_i)),$$

and

$$L(\theta, \alpha) = \int \ln(D_\alpha) p^* d\mu + \int \ln(1 - D_\alpha) p_\theta d\mu.$$

We will need the following regularity assumptions:

**Assumptions** ( $H_{\text{reg}}$ )

- ( $H_D$ ) There exists  $\kappa \in (0, 1/2)$  such that, for all  $\alpha \in \Lambda$ ,  $\kappa \leq D_\alpha \leq 1 - \kappa$ . In addition, the function  $(x, \alpha) \mapsto D_\alpha(x)$  is of class  $C^1$ , with a uniformly bounded differential.
- ( $H_G$ ) For all  $z \in \mathbb{R}^d$ , the function  $\theta \mapsto G_\theta(z)$  is of class  $C^1$ , uniformly bounded, with a uniformly bounded differential.
- ( $H_p$ ) For all  $x \in E$ , the function  $\theta \mapsto p_\theta(x)$  is of class  $C^1$ , uniformly bounded, with a uniformly bounded differential.

Note that under ( $H_D$ ), all discriminators in  $\{D_\alpha\}_{\alpha \in \Lambda}$  are  $\theta$ -admissible, whatever  $\theta$ . All of these requirements are classic regularity conditions for statistical models, which imply in particular that the functions  $\hat{L}(\theta, \alpha)$  and  $L(\theta, \alpha)$  are continuous. Therefore, the compactness of  $\Theta$  guarantees that  $\hat{\theta}$  and  $\bar{\theta}$  exist. Conditions for the existence of  $\theta^*$  are given in Theorem 2.2.

We have known since Theorem 3.1 that if the available class of discriminators  $\mathcal{D}$  approaches the optimal discriminator  $D_{\hat{\theta}}^*$  by a distance not more than  $\varepsilon$ , then  $D_{\text{JS}}(p^*, p_{\hat{\theta}}) - D_{\text{JS}}(p^*, p_{\theta^*}) = O(\varepsilon^2)$ . It is therefore reasonable to expect that, asymptotically, the difference  $D_{\text{JS}}(p^*, p_{\hat{\theta}}) - D_{\text{JS}}(p^*, p_{\theta^*})$  will not be larger than a term proportional to  $\varepsilon^2$ , in some probabilistic sense. This is precisely the result of Theorem 4.1 below. In fact, most articles to date have focused on the development and analysis of optimization procedures (typically, stochastic-gradient-type algorithms) to compute  $\hat{\theta}$ , without really questioning its convergence properties as the data set grows. Although our statistical results are theoretical in nature, we believe

that they are complementary to the optimization literature, insofar as they offer guarantees on the validity of the algorithms.

In addition to the regularity hypotheses, we will need the following requirement, which is a stronger version of  $(H_\varepsilon)$ :

**Assumption  $(H'_\varepsilon)$**  There exist  $\varepsilon > 0$  and  $m \in (0, 1/2)$  such that: for all  $\theta \in \Theta$ , there exists  $D \in \mathcal{D}$  such that  $m \leq D \leq 1 - m$  and  $\|D - D_\theta^*\|_2 \leq \varepsilon$ .

We are ready to state our first statistical theorem.

**THEOREM 4.1.** *Assume that, for some  $M > 0$ ,  $p^* \leq M$  and  $p_\theta \leq M$  for all  $\theta \in \Theta$ . Then, under Assumptions  $(H_{\text{reg}})$  and  $(H'_\varepsilon)$  with  $\varepsilon < 1/(2M)$ , one has*

$$\mathbb{E}D_{\text{JS}}(p^*, p_{\hat{\theta}}) - D_{\text{JS}}(p^*, p_{\theta^*}) = \mathcal{O}\left(\varepsilon^2 + \frac{1}{\sqrt{n}}\right).$$

**REMARK 4.1.** *The constant hidden in the  $\mathcal{O}$  term scales as  $p + q$ . Knowing that (deep) neural networks, and thus GANs, are often used in the so-called over-parameterized regime (i.e., when the number of parameters exceeds the number of examples), this limits the impact of the result in the neural network context, at least when  $p + q$  is large with respect to  $\sqrt{n}$ . For instance, successful applications of GANs on common datasets such as LSUN ( $\sqrt{n} \approx 1740$ ) and FACES ( $\sqrt{n} \approx 590$ ) make use of more than 1 500 000 parameters (Radford et al., 2016).*

**PROOF.** Fix  $\varepsilon \in (0, 1/(2M))$  as in Assumption  $(H'_\varepsilon)$ , and choose  $\hat{D} \in \mathcal{D}$  such that  $m \leq \hat{D} \leq 1 - m$  and  $\|\hat{D} - D_\theta^*\|_2 \leq \varepsilon$ . By repeating the arguments of the proof of Theorem 3.1 (with  $\hat{\theta}$  instead of  $\bar{\theta}$ ), we conclude that there exists a constant  $c_1 > 0$  such that

$$2D_{\text{JS}}(p^*, p_{\hat{\theta}}) \leq c_1 \varepsilon^2 + L(\hat{\theta}, \hat{D}) + \ln 4 \leq c_1 \varepsilon^2 + \sup_{\alpha \in \Lambda} L(\hat{\theta}, \alpha) + \ln 4.$$

Therefore,

$$\begin{aligned} 2D_{\text{JS}}(p^*, p_{\hat{\theta}}) &\leq c_1 \varepsilon^2 + \sup_{\theta \in \Theta, \alpha \in \Lambda} |\hat{L}(\theta, \alpha) - L(\theta, \alpha)| + \sup_{\alpha \in \Lambda} \hat{L}(\hat{\theta}, \alpha) + \ln 4 \\ &= c_1 \varepsilon^2 + \sup_{\theta \in \Theta, \alpha \in \Lambda} |\hat{L}(\theta, \alpha) - L(\theta, \alpha)| + \inf_{\theta \in \Theta} \sup_{\alpha \in \Lambda} \hat{L}(\theta, \alpha) + \ln 4 \\ &\quad \text{(by definition of } \hat{\theta}\text{)} \\ &\leq c_1 \varepsilon^2 + 2 \sup_{\theta \in \Theta, \alpha \in \Lambda} |\hat{L}(\theta, \alpha) - L(\theta, \alpha)| + \inf_{\theta \in \Theta} \sup_{\alpha \in \Lambda} L(\theta, \alpha) + \ln 4. \end{aligned}$$

So,

$$\begin{aligned}
2D_{\text{JS}}(p^*, p_{\hat{\theta}}) &\leq c_1 \varepsilon^2 + 2 \sup_{\theta \in \Theta, \alpha \in \Lambda} |\hat{L}(\theta, \alpha) - L(\theta, \alpha)| + \inf_{\theta \in \Theta} \sup_{D \in \mathcal{D}_\infty} L(\theta, D) + \ln 4 \\
&= c_1 \varepsilon^2 + 2 \sup_{\theta \in \Theta, \alpha \in \Lambda} |\hat{L}(\theta, \alpha) - L(\theta, \alpha)| + L(\theta^*, D_{\theta^*}^*) + \ln 4 \\
&\quad \text{(by definition of } \theta^*) \\
&= c_1 \varepsilon^2 + 2D_{\text{JS}}(p^*, p_{\theta^*}) + 2 \sup_{\theta \in \Theta, \alpha \in \Lambda} |\hat{L}(\theta, \alpha) - L(\theta, \alpha)|.
\end{aligned}$$

Thus, letting  $c_2 = c_1/2$ , we have

$$(4.1) \quad D_{\text{JS}}(p^*, p_{\hat{\theta}}) - D_{\text{JS}}(p^*, p_{\theta^*}) \leq c_2 \varepsilon^2 + \sup_{\theta \in \Theta, \alpha \in \Lambda} |\hat{L}(\theta, \alpha) - L(\theta, \alpha)|.$$

Clearly, under Assumptions  $(H_D)$ ,  $(H_G)$ , and  $(H_p)$ ,  $(\hat{L}(\theta, \alpha) - L(\theta, \alpha))_{\theta \in \Theta, \alpha \in \Lambda}$  is a separable subgaussian process (e.g., [van Handel, 2016](#), Chapter 5) for the distance  $d = S \|\cdot\|/\sqrt{n}$ , where  $\|\cdot\|$  is the standard Euclidean norm on  $\mathbb{R}^p \times \mathbb{R}^q$  and  $S > 0$  depends only on the bounds in  $(H_D)$  and  $(H_G)$ . Let  $N(\Theta \times \Lambda, \|\cdot\|, u)$  denote the  $u$ -covering number of  $\Theta \times \Lambda$  for the distance  $\|\cdot\|$ . Then, by Dudley's inequality ([van Handel, 2016](#), Corollary 5.25),

$$(4.2) \quad \mathbb{E} \sup_{\theta \in \Theta, \alpha \in \Lambda} |\hat{L}(\theta, \alpha) - L(\theta, \alpha)| \leq \frac{12S}{\sqrt{n}} \int_0^\infty \sqrt{\ln(N(\Theta \times \Lambda, \|\cdot\|, u))} du.$$

Since  $\Theta$  and  $\Lambda$  are bounded, there exists  $r > 0$  such that  $N(\Theta \times \Lambda, \|\cdot\|, u) = 1$  for  $u \geq r$  and

$$N(\Theta \times \Lambda, \|\cdot\|, u) = \mathcal{O}\left(\left(\frac{\sqrt{p+q}}{u}\right)^{p+q}\right) \quad \text{for } u < r.$$

Combining this inequality with (4.1) and (4.2), we obtain

$$\mathbb{E}D_{\text{JS}}(p^*, p_{\hat{\theta}}) - D_{\text{JS}}(p^*, p_{\theta^*}) \leq c_3 \left( \varepsilon^2 + \frac{1}{\sqrt{n}} \right),$$

for some positive constant  $c_3$  that scales as  $p + q$ . The conclusion follows by observing that, by (2.2),

$$D_{\text{JS}}(p^*, p_{\theta^*}) \leq D_{\text{JS}}(p^*, p_{\hat{\theta}}).$$

□

Theorem 4.1 is illustrated in Figure 1, which shows the approximate values of  $\mathbb{E}D_{\text{JS}}(p^*, p_{\hat{\theta}})$ . We took  $p^*(x) = \frac{e^{-x/s}}{s(1+e^{-x/s})^2}$  (centered logistic density with scale parameter  $s = 0.33$ ), and let  $\mathcal{G}$  and  $\mathcal{D}$  be two fully connected neural networks parameterized by weights and offsets. The noise random variable  $Z$  follows a uniform

distribution on  $[0, 1]$ , and the parameters of  $\mathcal{G}$  and  $\mathcal{D}$  are chosen in a sufficiently large compact set. In order to illustrate the impact of  $\varepsilon$  in Theorem 4.1, we fixed the sample size to a large  $n = 100\,000$  and varied the number of layers of the discriminators from 2 to 5, keeping in mind that a larger number of layers results in a smaller  $\varepsilon$ . To diversify the setting, we also varied the number of layers of the generators from 2 to 3. The expectation  $\mathbb{E}D_{\text{JS}}(p^*, p_{\hat{\theta}})$  was estimated by averaging over 30 repetitions (the number of runs has been reduced for time complexity limitations). Note that we do not pay attention to the exact value of the constant term  $D_{\text{JS}}(p^*, p_{\theta^*})$ , which is intractable in our setting.

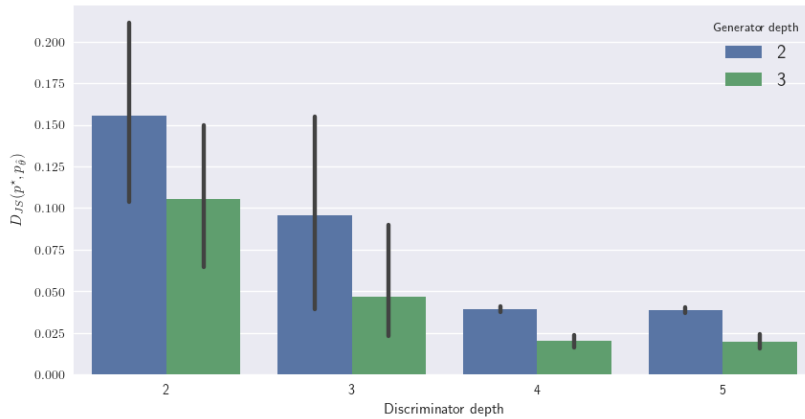
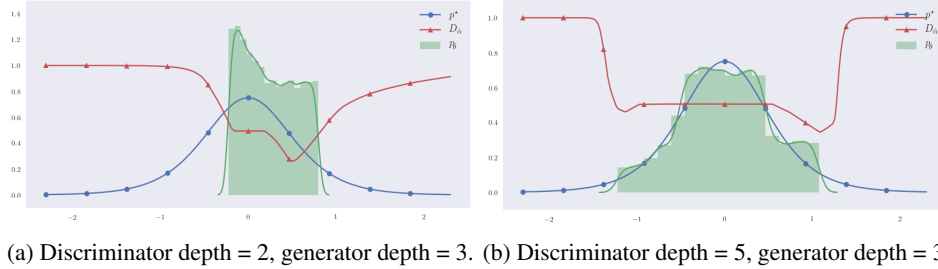


Fig 1: Bar plots of the Jensen-Shannon divergence  $D_{\text{JS}}(p^*, p_{\hat{\theta}})$  with respect to the number of layers (depth) of both the discriminators and generators. The height of each rectangle estimates  $\mathbb{E}D_{\text{JS}}(p^*, p_{\hat{\theta}})$ .

Figure 1 highlights that  $\mathbb{E}D_{\text{JS}}(p^*, p_{\hat{\theta}})$  approaches the value  $D_{\text{JS}}(p^*, p_{\theta^*})$  as  $\varepsilon \downarrow 0$ , i.e., as the discriminator depth increases, given that the contribution of  $1/\sqrt{n}$  is certainly negligible for  $n = 100\,000$ . Figure 2 shows the target density  $p^*$  vs. the histograms and kernel estimates of 100 000 data sampled from  $G_{\hat{\theta}}(Z)$ , in the two cases: (discriminator depth = 2, generator depth = 3) and (discriminator depth = 5, generator depth = 3). In accordance with the decrease of  $\mathbb{E}D_{\text{JS}}(p^*, p_{\hat{\theta}})$ , the estimation of the true distribution  $p^*$  improves when  $\varepsilon$  becomes small.

*Some comments on the optimization scheme.* Numerical optimization is quite a tough point for GANs, partly due to nonconvex-concavity of the saddle point problem described in equation (1.2) and the nondifferentiability of the objective function. This motivates a very active line of research (e.g., Goodfellow et al., 2014; Nowozin et al., 2016; Arjovsky et al., 2017; Arjovsky and Bottou, 2017), which



(a) Discriminator depth = 2, generator depth = 3. (b) Discriminator depth = 5, generator depth = 3.

Fig 2: True density  $p^*$ , histograms, and kernel estimates (continuous line) of 100000 data sampled from  $G_{\hat{\theta}}(Z)$ . Also shown is the final discriminator  $D_{\hat{\alpha}}$ .

aims at transforming the objective into a more convenient function and devising efficient algorithms. In the present paper, since we are interested in original GANs, the algorithmic approach described by Goodfellow et al. (2014) is adopted, and numerical optimization is performed thanks to the machine learning framework TensorFlow (Abadi et al., 2015), working with gradient descent based on automatic differentiation. As proposed by Goodfellow et al. (2014), the objective function  $\theta \mapsto \sup_{\alpha \in \Lambda} \hat{L}(\theta, \alpha)$  is not directly minimized. We used instead an alternated procedure, which consists in iterating (a few hundred times in our examples) the following two steps:

- (i) For a fixed value of  $\theta$  and from a given value of  $\alpha$ , perform 10 ascent steps on  $\hat{L}(\theta, \cdot)$ ;
- (ii) For a fixed value of  $\alpha$  and from a given value of  $\theta$ , perform 1 descent step on  $\theta \mapsto -\sum_{i=1}^n \ln(D_{\alpha} \circ G_{\theta}(Z_i))$  (instead of  $\theta \mapsto \sum_{i=1}^n \ln(1 - D_{\alpha} \circ G_{\theta}(Z_i))$ ).

This alternated procedure is motivated by two reasons. First, for a given  $\theta$ , approximating  $\sup_{\alpha \in \Lambda} \hat{L}(\theta, \alpha)$  is computationally prohibitive and may result in overfitting the finite training sample (Goodfellow et al., 2014). This can be explained by the shape of the function  $\theta \mapsto \sup_{\alpha \in \Lambda} \hat{L}(\theta, \alpha)$ , which may be almost piecewise constant, resulting in a zero gradient almost everywhere (or at best very low; see Arjovsky et al., 2017). Next, empirically,  $-\ln(D_{\alpha} \circ G_{\theta}(Z_i))$  provides bigger gradients than  $\ln(1 - D_{\alpha} \circ G_{\theta}(Z_i))$ , resulting in a more powerful algorithm than the original version, while leading to the same minimizers.

In all our experiments, the learning rates needed in gradient steps were fixed and tuned by hand, in order to prevent divergence. In addition, since our main objective is to focus on illustrating the statistical properties of GANs rather than delving into optimization issues, we decided to perform mini-batch gradient updates instead of stochastic ones (that is, new observations of  $X$  and  $Z$  are not sampled at each step

of the procedure). This is different from what is done in the original algorithm of [Goodfellow et al. \(2014\)](#).

All in all, we realize that our numerical approach—although widely adopted by the machine learning community—may fail to locate the desired estimator  $\hat{\theta}$  (i.e., the exact minimizer in  $\theta$  of  $\sup_{\alpha \in \Lambda} \hat{L}(\theta, \alpha)$ ) in more complex contexts than those presented in the present paper. It is nevertheless sufficient for our objective, which is limited to illustrating the theoretical results with a few simple examples.

*4.2. Asymptotic properties of  $\hat{\theta}$ .* Theorem 4.1 states a result relative to the excess of Jensen-Shannon error  $D_{\text{JS}}(p^*, p_{\hat{\theta}}) - D_{\text{JS}}(p^*, p_{\theta^*})$ . We now examine the convergence properties of the parameter  $\hat{\theta}$  itself as the sample size  $n$  grows. We would typically like to find reasonable conditions ensuring that  $\hat{\theta} \rightarrow \bar{\theta}$  almost surely as  $n \rightarrow \infty$ . To reach this goal, we first need to strengthen a bit the Assumptions  $(H_{\text{reg}})$ , as follows:

**Assumptions  $(H'_{\text{reg}})$**

- $(H'_D)$  There exists  $\kappa \in (0, 1/2)$  such that, for all  $\alpha \in \Lambda$ ,  $\kappa \leq D_\alpha \leq 1 - \kappa$ . In addition, the function  $(x, \alpha) \mapsto D_\alpha(x)$  is of class  $C^2$ , with differentials of order 1 and 2 uniformly bounded.
- $(H'_G)$  For all  $z \in \mathbb{R}^{d'}$ , the function  $\theta \mapsto G_\theta(z)$  is of class  $C^2$ , uniformly bounded, with differentials of order 1 and 2 uniformly bounded.
- $(H'_p)$  For all  $x \in E$ , the function  $\theta \mapsto p_\theta(x)$  is of class  $C^2$ , uniformly bounded, with differentials of order 1 and 2 uniformly bounded.

It is easy to verify that under these assumptions the partial functions  $\theta \mapsto \hat{L}(\theta, \alpha)$  (respectively,  $\theta \mapsto L(\theta, \alpha)$ ) and  $\alpha \mapsto \hat{L}(\theta, \alpha)$  (respectively,  $\alpha \mapsto L(\theta, \alpha)$ ) are of class  $C^2$ . Throughout, we let  $\theta = (\theta_1, \dots, \theta_p)$ ,  $\alpha = (\alpha_1, \dots, \alpha_q)$ , and denote by  $\frac{\partial}{\partial \theta_i}$  and  $\frac{\partial}{\partial \alpha_j}$  the partial derivative operations with respect to  $\theta_i$  and  $\alpha_j$ . The next lemma will be of constant utility. In order not to burden the text, its proof is given in Section 6.

LEMMA 4.1. *Under Assumptions  $(H'_{\text{reg}})$ ,  $\forall (a, b, c, d) \in \{0, 1, 2\}^4$  such that  $a + b \leq 2$  and  $c + d \leq 2$ , one has*

$$\sup_{\theta \in \Theta, \alpha \in \Lambda} \left| \frac{\partial^{a+b+c+d}}{\partial \theta_i^a \partial \theta_j^b \partial \alpha_\ell^c \partial \alpha_m^d} \hat{L}(\theta, \alpha) - \frac{\partial^{a+b+c+d}}{\partial \theta_i^a \partial \theta_j^b \partial \alpha_\ell^c \partial \alpha_m^d} L(\theta, \alpha) \right| \rightarrow 0$$

*almost surely, for all  $(i, j) \in \{1, \dots, p\}^2$  and  $(\ell, m) \in \{1, \dots, q\}^2$ .*

We recall that  $\bar{\theta} \in \Theta$  is such that

$$\sup_{\alpha \in \Lambda} L(\bar{\theta}, \alpha) \leq \sup_{\alpha \in \Lambda} L(\theta, \alpha), \quad \forall \theta \in \Theta,$$



and insist that  $\bar{\theta}$  exists under  $(H'_{\text{reg}})$  by continuity of the map  $\theta \mapsto \sup_{\alpha \in \Lambda} L(\theta, \alpha)$ . Similarly, there exists  $\bar{\alpha} \in \Lambda$  such that

$$L(\bar{\theta}, \bar{\alpha}) \geq L(\bar{\theta}, \alpha), \quad \forall \alpha \in \Lambda.$$

The following assumption ensures that  $\bar{\theta}$  and  $\bar{\alpha}$  are uniquely defined, which is of course a key hypothesis for our estimation objective. Throughout, the notation  $S^\circ$  (respectively,  $\partial S$ ) stands for the interior (respectively, the boundary) of the set  $S$ .

**Assumption  $(H_1)$**  The pair  $(\bar{\theta}, \bar{\alpha})$  is unique and belongs to  $\Theta^\circ \times \Lambda^\circ$ .

Finally, in addition to  $\hat{\theta}$ , we let  $\hat{\alpha} \in \Lambda$  be such that

$$\hat{L}(\hat{\theta}, \hat{\alpha}) \geq \hat{L}(\hat{\theta}, \alpha), \quad \forall \alpha \in \Lambda.$$

**THEOREM 4.2.** *Under Assumptions  $(H'_{\text{reg}})$  and  $(H_1)$ , one has*

$$\hat{\theta} \rightarrow \bar{\theta} \quad \text{almost surely} \quad \text{and} \quad \hat{\alpha} \rightarrow \bar{\alpha} \quad \text{almost surely.}$$

**PROOF.** We write

$$\begin{aligned} & \left| \sup_{\alpha \in \Lambda} L(\hat{\theta}, \alpha) - \sup_{\alpha \in \Lambda} L(\bar{\theta}, \alpha) \right| \\ & \leq \left| \sup_{\alpha \in \Lambda} L(\hat{\theta}, \alpha) - \sup_{\alpha \in \Lambda} \hat{L}(\hat{\theta}, \alpha) \right| + \left| \inf_{\theta \in \Theta} \sup_{\alpha \in \Lambda} \hat{L}(\theta, \alpha) - \inf_{\theta \in \Theta} \sup_{\alpha \in \Lambda} L(\theta, \alpha) \right| \\ & \leq 2 \sup_{\theta \in \Theta, \alpha \in \Lambda} |\hat{L}(\theta, \alpha) - L(\theta, \alpha)|. \end{aligned}$$

Thus, by Lemma 4.1,  $\sup_{\alpha \in \Lambda} L(\hat{\theta}, \alpha) \rightarrow \sup_{\alpha \in \Lambda} L(\bar{\theta}, \alpha)$  almost surely. In the lines that follow, we make more transparent the dependence of  $\hat{\theta}$  in the sample size  $n$  and set  $\hat{\theta}_n \stackrel{\text{def}}{=} \hat{\theta}$ . Since  $\hat{\theta}_n \in \Theta$  and  $\Theta$  is compact, we can extract from any subsequence of  $(\hat{\theta}_n)_n$  a subsequence  $(\hat{\theta}_{n_k})_k$  such that  $\hat{\theta}_{n_k} \rightarrow z \in \Theta$  (with  $n_k = n_k(\omega)$ , i.e., it is almost surely defined). By continuity of the function  $\theta \mapsto \sup_{\alpha \in \Lambda} L(\theta, \alpha)$ , we deduce that  $\sup_{\alpha \in \Lambda} \hat{L}(\hat{\theta}_{n_k}, \alpha) \rightarrow \sup_{\alpha \in \Lambda} L(z, \alpha)$ , and so  $\sup_{\alpha \in \Lambda} L(z, \alpha) = \sup_{\alpha \in \Lambda} L(\bar{\theta}, \alpha)$ . Since  $\bar{\theta}$  is unique by  $(H_1)$ , we have  $z = \bar{\theta}$ . In conclusion, we can extract from each subsequence of  $(\hat{\theta}_n)_n$  a subsequence that converges towards  $\bar{\theta}$ : this shows that  $\hat{\theta}_n \rightarrow \bar{\theta}$  almost surely.

Finally, we have

$$\begin{aligned} & |L(\bar{\theta}, \hat{\alpha}) - L(\bar{\theta}, \bar{\alpha})| \\ & \leq |L(\bar{\theta}, \hat{\alpha}) - L(\hat{\theta}, \hat{\alpha})| + |L(\hat{\theta}, \hat{\alpha}) - \hat{L}(\hat{\theta}, \hat{\alpha})| + |\hat{L}(\hat{\theta}, \hat{\alpha}) - L(\bar{\theta}, \bar{\alpha})| \\ & = |L(\bar{\theta}, \hat{\alpha}) - L(\hat{\theta}, \hat{\alpha})| + |L(\hat{\theta}, \hat{\alpha}) - \hat{L}(\hat{\theta}, \hat{\alpha})| \\ & \quad + \left| \inf_{\theta \in \Theta} \sup_{\alpha \in \Lambda} \hat{L}(\theta, \alpha) - \inf_{\theta \in \Theta} \sup_{\alpha \in \Lambda} L(\theta, \alpha) \right| \\ & \leq \sup_{\alpha \in \Lambda} |L(\bar{\theta}, \alpha) - L(\hat{\theta}, \alpha)| + 2 \sup_{\theta \in \Theta, \alpha \in \Lambda} |\hat{L}(\theta, \alpha) - L(\theta, \alpha)|. \end{aligned}$$

Using Assumptions  $(H'_D)$  and  $(H'_p)$ , and the fact that  $\hat{\theta} \rightarrow \bar{\theta}$  almost surely, we see that the first term above tends to zero. The second one vanishes asymptotically by Lemma 4.1, and we conclude that  $L(\bar{\theta}, \hat{\alpha}) \rightarrow L(\bar{\theta}, \bar{\alpha})$  almost surely. Since  $\hat{\alpha} \in \Lambda$  and  $\Lambda$  is compact, we may argue as in the first part of the proof and deduce from the uniqueness of  $\bar{\alpha}$  that  $\hat{\alpha} \rightarrow \bar{\alpha}$  almost surely.  $\square$

To illustrate the result of Theorem 4.2, we undertook a series of small numerical experiments with three choices for the triplet (true  $p^*$  + generator model  $\mathcal{P} = \{p_\theta\}_{\theta \in \Theta}$  + discriminator family  $\mathcal{D} = \{D_\alpha\}_{\alpha \in \Lambda}$ ), which we respectively call the **Laplace-Gaussian**, **Claw-Gaussian**, and **Exponential-Uniform** model. They are summarized in Table 1. We are aware that more elaborate models (involving, for example, neural networks) can be designed and implemented. However, our objective is not to conduct a series of extensive simulations, but simply to illustrate our theoretical results with a few graphs to get some better intuition and provide a sanity check. We stress in particular that these experiments are in one dimension and are therefore very limited compared to the way GANs algorithms are typically used in practice.

Model	$p^*$	$\mathcal{P} = \{p_\theta\}_{\theta \in \Theta}$	$\mathcal{D} = \{D_\alpha\}_{\alpha \in \Lambda}$
<b>Laplace-Gaussian</b>	$\frac{1}{2b} e^{-\frac{ x }{b}}$ $b = 1.5$	$\frac{1}{\sqrt{2\pi}\theta} e^{-\frac{x^2}{2\theta^2}}$ $\Theta = [10^{-1}, 10^3]$	$\frac{1}{1 + \frac{\alpha_1}{\alpha_0} e^{\frac{\alpha_1^2}{2}(\alpha_1^{-2} - \alpha_0^{-2})}}$ $\Lambda = \Theta \times \Theta$
<b>Claw-Gaussian</b>	$p_{\text{claw}}(x)$	$\frac{1}{\sqrt{2\pi}\theta} e^{-\frac{x^2}{2\theta^2}}$ $\Theta = [10^{-1}, 10^3]$	$\frac{1}{1 + \frac{\alpha_1}{\alpha_0} e^{\frac{\alpha_1^2}{2}(\alpha_1^{-2} - \alpha_0^{-2})}}$ $\Lambda = \Theta \times \Theta$
<b>Exponential-Uniform</b>	$\lambda e^{-\lambda x}$ $\lambda = 1$	$\frac{1}{\theta} \mathbf{1}_{[0, \theta]}(x)$ $\Theta = [10^{-3}, 10^3]$	$\frac{1}{1 + \frac{\alpha_1}{\alpha_0} e^{\frac{\alpha_1^2}{2}(\alpha_1^{-2} - \alpha_0^{-2})}}$ $\Lambda = \Theta \times \Theta$

TABLE 1  
Triplets used in the numerical experiments.

Figure 3 shows the densities  $p^*$ . We recall that the claw density on  $(-\infty, \infty)$  takes the form

$$p_{\text{claw}} = \frac{1}{2} \varphi(0, 1) + \frac{1}{10} (\varphi(-1, 0.1) + \varphi(-0.5, 0.1) + \varphi(0, 0.1) + \varphi(0.5, 0.1) + \varphi(1, 0.1)),$$

where  $\varphi(\mu, \sigma)$  is a Gaussian density with mean  $\mu$  and standard deviation  $\sigma$  (this density is borrowed from Devroye, 1997).

In the **Laplace-Gaussian** and **Claw-Gaussian** examples, the densities  $p_\theta$  are centered Gaussian, parameterized by their standard deviation parameter  $\theta$ . The random variable  $Z$  is uniform on  $[0, 1]$  and the natural family of generators associated

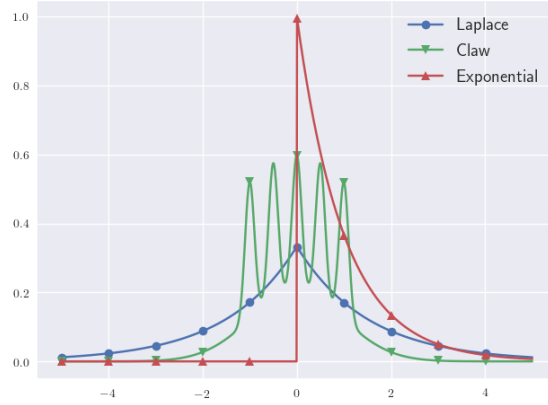


Fig 3: Probability density functions  $p^*$  used in the numerical experiments.

with the model  $\mathcal{P} = \{p_\theta\}_{\theta \in \Theta}$  is  $\mathcal{G} = \{G_\theta\}_{\theta \in \Theta}$ , where each  $G_\theta$  is the generalized inverse of the cumulative distribution function of  $p_\theta$  (because  $G_\theta(Z) \stackrel{\mathcal{L}}{=} p_\theta d\mu$ ). The rationale behind our choice for the discriminators is based on the form of the optimal discriminator  $D_\theta^*$  described in (2.1): starting from

$$D_\theta^* = \frac{p^*}{p^* + p_\theta}, \quad \theta \in \Theta,$$

we logically consider the following ratio

$$D_\alpha = \frac{p_{\alpha_1}}{p_{\alpha_1} + p_{\alpha_0}}, \quad \alpha = (\alpha_0, \alpha_1) \in \Lambda = \Theta \times \Theta.$$

Figure 4 (**Laplace-Gaussian**), Figure 5 (**Claw-Gaussian**), and Figure 6 (**Exponential-Uniform**) show the boxplots of the differences  $\hat{\theta} - \bar{\theta}$  over 200 repetitions, for a sample size  $n$  varying from 10 to 10000. In these experiments, the parameter  $\bar{\theta}$  is obtained by averaging the  $\hat{\theta}$  for the largest sample size  $n$ . In accordance with Theorem 4.2, the size of the boxplots shrinks around 0 when  $n$  increases, thus showing that the estimated parameter  $\hat{\theta}$  is getting closer and closer to  $\bar{\theta}$ . Before analyzing at which rate this convergence occurs, we may have a look at Figure 7, which plots the estimated density  $p_{\hat{\theta}}$  (for  $n = 10000$ ) vs. the true density  $p^*$ . It also shows the discriminator  $D_{\hat{\alpha}}$ , together with the initial density  $p_{\theta_{\text{init}}}$  and the initial discriminator  $D_{\alpha_{\text{init}}}$  fed into the optimization algorithm. We note that in the three models,  $D_{\hat{\alpha}}$  is almost identically  $1/2$ , meaning that it is impossible to discriminate between the original observations and those generated by  $p_{\hat{\theta}}$ .

In line with the above, our next step is to state a central limit theorem for  $\hat{\theta}$ . Although simple to understand, this result requires additional assumptions and some

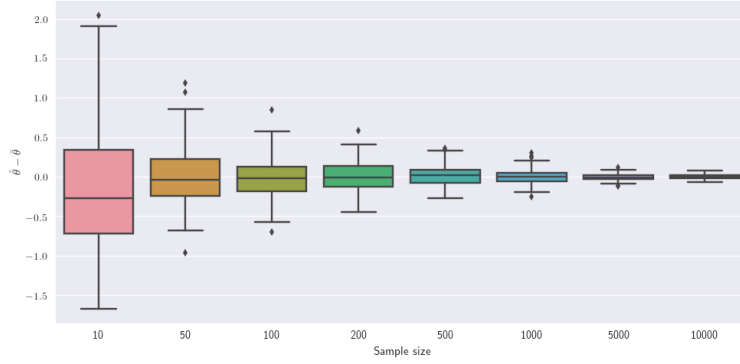


Fig 4: Boxplots of  $\hat{\theta} - \bar{\theta}$  for different sample sizes (**Laplace-Gaussian** model, 200 repetitions).

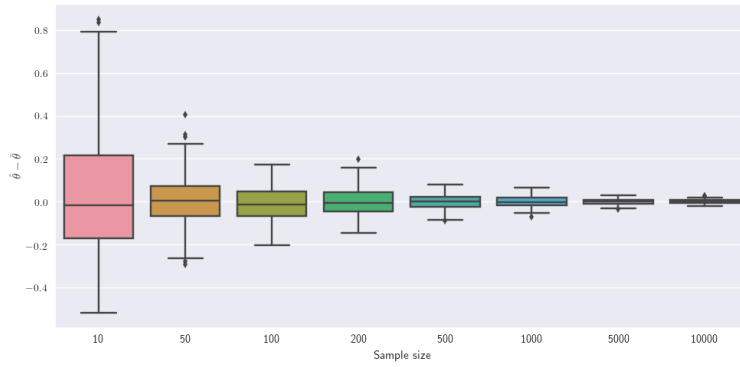


Fig 5: Boxplots of  $\hat{\theta} - \bar{\theta}$  for different sample sizes (**Claw-Gaussian** model, 200 repetitions).

technical prerequisites. One first needs to ensure that the function  $(\theta, \alpha) \mapsto L(\theta, \alpha)$  is regular enough in a neighborhood of  $(\bar{\theta}, \bar{\alpha})$ . This is captured by the following set of assumptions, which require in particular the uniqueness of the maximizer of the function  $\alpha \mapsto L(\theta, \alpha)$  for a  $\theta$  around  $\bar{\theta}$ . For a function  $F : \Theta \rightarrow \mathbb{R}$  (respectively,  $G : \Theta \times \Lambda \rightarrow \mathbb{R}$ ), we let  $HF(\theta)$  (respectively,  $H_1G(\theta, \alpha)$  and  $H_2G(\theta, \alpha)$ ) be the Hessian matrix of the function  $\theta \mapsto F(\theta)$  (respectively,  $\theta \mapsto G(\theta, \alpha)$  and  $\alpha \mapsto G(\theta, \alpha)$ ) computed at  $\theta$  (respectively, at  $\theta$  and  $\alpha$ ).

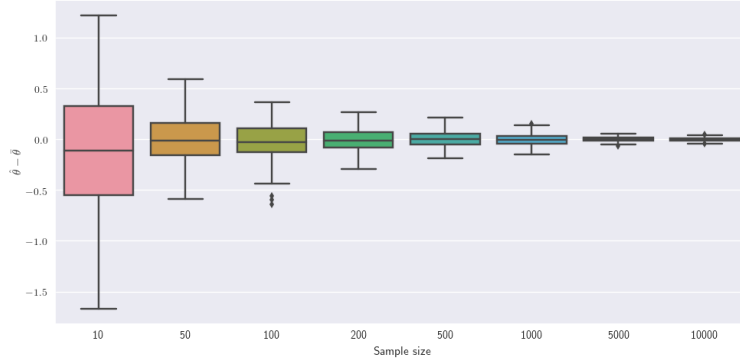


Fig 6: Boxplots of  $\hat{\theta} - \bar{\theta}$  for different sample sizes (**Exponential-Uniform** model, 200 repetitions).

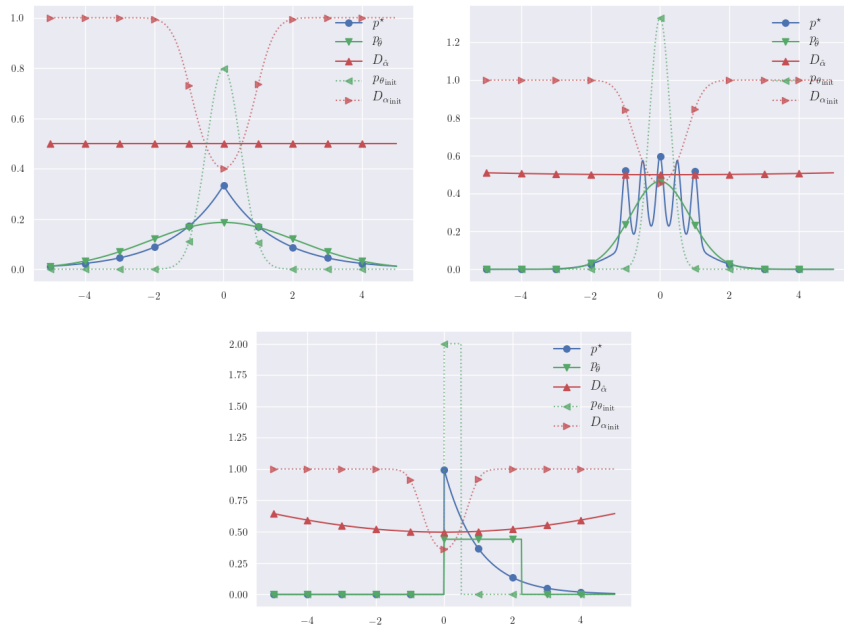


Fig 7: True density  $p^*$ , estimated density  $p_{\hat{\theta}}$ , and discriminator  $D_{\hat{\alpha}}$  for  $n = 10000$  (from left to right: **Laplace-Gaussian**, **Claw-Gaussian**, and **Exponential-Uniform** model). Also shown are the initial density  $p_{\theta_{\text{init}}}$  and the initial discriminator  $D_{\alpha_{\text{init}}}$  fed into the optimization algorithm.

**Assumptions** ( $H_{\text{loc}}$ )

( $H_U$ ) There exists a neighborhood  $U$  of  $\bar{\theta}$  and a function  $\alpha : U \rightarrow \Lambda$  such that

$$\arg \max_{\alpha \in \Lambda} L(\theta, \alpha) = \{\alpha(\theta)\}, \quad \forall \theta \in U.$$

( $H_V$ ) The Hessian matrix  $HV(\bar{\theta})$  is invertible, where  $V(\theta) \stackrel{\text{def}}{=} L(\theta, \alpha(\theta))$ .

( $H_H$ ) The Hessian matrix  $H_2L(\bar{\theta}, \bar{\alpha})$  is invertible.

We stress that under Assumption ( $H_U$ ), there is for each  $\theta \in U$  a unique  $\alpha(\theta) \in \Lambda$  such that  $L(\theta, \alpha(\theta)) = \sup_{\alpha \in \Lambda} L(\theta, \alpha)$ . We also note that  $\alpha(\bar{\theta}) = \bar{\alpha}$  under ( $H_1$ ). We still need some notation before we state the central limit theorem. For a function  $f(\theta, \alpha)$ ,  $\nabla_1 f(\theta, \alpha)$  (respectively,  $\nabla_2 f(\theta, \alpha)$ ) means the gradient of the function  $\theta \mapsto f(\theta, \alpha)$  (respectively, the function  $\alpha \mapsto f(\theta, \alpha)$ ) computed at  $\theta$  (respectively, at  $\alpha$ ). For a function  $g(t)$ ,  $J(g)_t$  is the Jacobian matrix of  $g$  computed at  $t$ . Observe that by the envelope theorem,

$$HV(\bar{\theta}) = H_1L(\bar{\theta}, \bar{\alpha}) + J(\nabla_1L(\bar{\theta}, \cdot))_{\bar{\alpha}}J(\alpha)_{\bar{\theta}},$$

where, by the chain rule,

$$J(\alpha)_{\bar{\theta}} = -H_2L(\bar{\theta}, \bar{\alpha})^{-1}J(\nabla_2L(\cdot, \bar{\alpha}))_{\bar{\theta}}.$$

Therefore, in Assumption ( $H_V$ ), the Hessian matrix  $HV(\bar{\theta})$  can be computed with the sole knowledge of  $L$ . Finally, we let

$$\ell_1(\theta, \alpha) = \ln D_\alpha(X_1) + \ln(1 - D_\alpha \circ G_\theta(Z_1)),$$

and denote by  $\xrightarrow{\mathcal{L}}$  the convergence in distribution.

**THEOREM 4.3.** *Under Assumptions ( $H'_{\text{reg}}$ ), ( $H_1$ ), and ( $H_{\text{loc}}$ ), one has*

$$\sqrt{n}(\hat{\theta} - \bar{\theta}) \xrightarrow{\mathcal{L}} Z,$$

where  $Z$  is a Gaussian random variable with mean 0 and covariance matrix

$$\mathbf{V} = \text{Var}[-HV(\bar{\theta})^{-1}\nabla_1\ell_1(\bar{\theta}, \bar{\alpha}) + HV(\bar{\theta})^{-1}J(\nabla_1L(\bar{\theta}, \cdot))_{\bar{\alpha}}H_2L(\bar{\theta}, \bar{\alpha})^{-1}\nabla_2\ell_1(\bar{\theta}, \bar{\alpha})].$$

The expression of the covariance is relatively complex and, unfortunately, cannot be simplified, even for a dimension of the parameter equal to 1. We note however that if  $Y$  is a random vector of  $\mathbb{R}^p$  whose components are bounded in absolute value by some  $\delta > 0$ , then the Euclidean norm of the covariance matrix of  $Y$  is bounded by  $4p\delta^2$ . But each component of the random vector of  $\mathbb{R}^p$  involved in

the covariance matrix  $\mathbf{V}$  is bounded in absolute value by  $Cpq^2$ , for some positive constant  $C$  resulting from Assumption  $(H'_{\text{reg}})$ . We conclude that the Euclidean norm of  $\mathbf{V}$  is bounded by  $4C^2p^3q^4$ . Thus, our statistical approach reveals that in the overparameterized regime (i.e, when  $p$  and  $q$  are very large compared to  $n$ ), the estimator  $\hat{\theta}$  has a large dispersion around  $\bar{\theta}$ , which may affects the performance of the algorithm.

Nevertheless, the take-home message of Theorem 4.3 is that the estimator  $\hat{\theta}$  is asymptotically normal, with a convergence rate of  $\sqrt{n}$ . This is illustrated in Figures 8, 9, and 10, which respectively show the histograms and kernel estimates of the distribution of  $\sqrt{n}(\hat{\theta} - \bar{\theta})$  for the **Laplace-Gaussian**, the **Claw-Gaussian**, and the **Exponential-Uniform** model in function of the sample size  $n$  (200 repetitions).

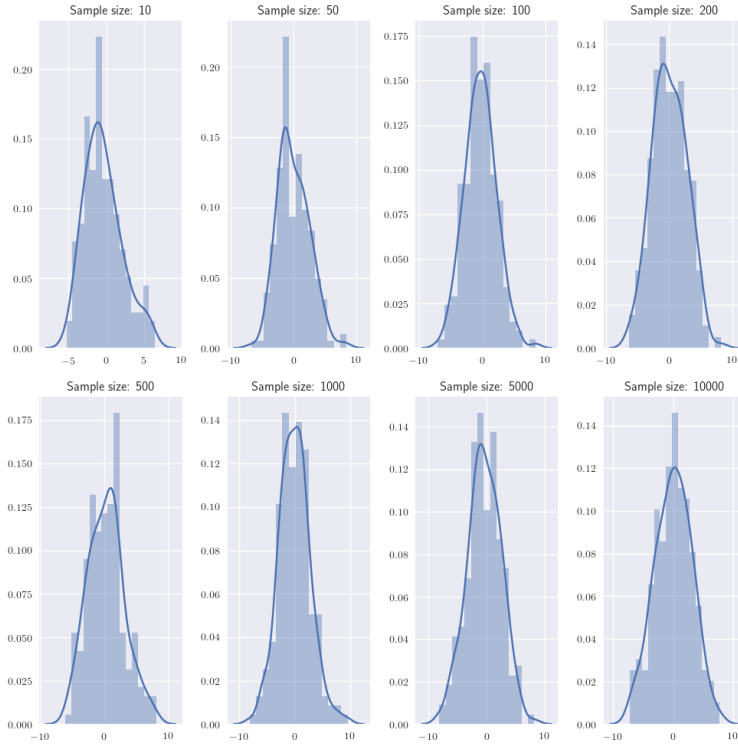


Fig 8: Histograms and kernel estimates (continuous line) of the distribution of  $\sqrt{n}(\hat{\theta} - \bar{\theta})$  for different sample sizes  $n$  (**Laplace-Gaussian** model, 200 repetitions).

PROOF. By technical Lemma 6.1, we can find under Assumptions  $(H'_{\text{reg}})$  and  $(H_1)$  an open set  $V \subset U \subset \Theta^\circ$  containing  $\bar{\theta}$  such that, for all  $\theta \in V$ ,  $\alpha(\theta) \in \Lambda^\circ$ .

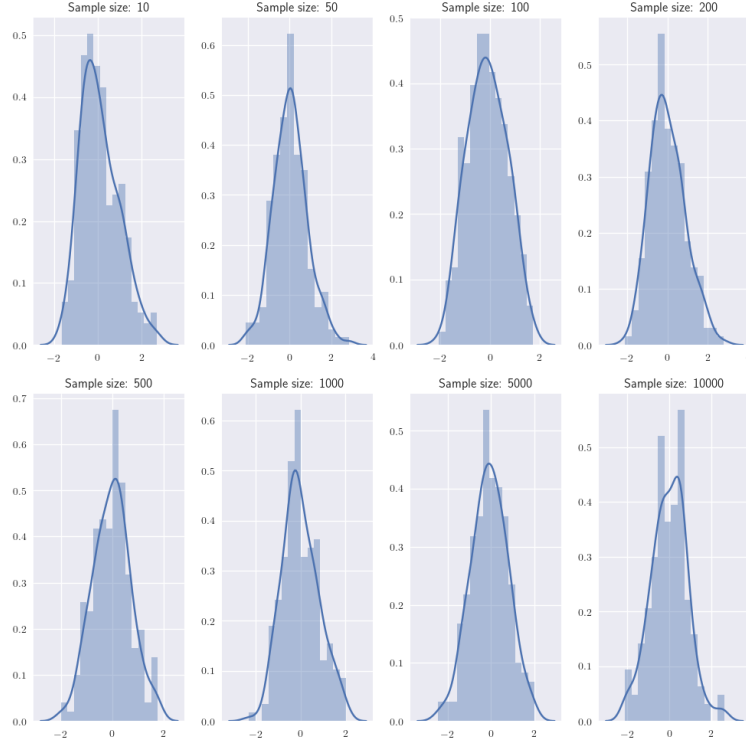


Fig 9: Histograms and kernel estimates (continuous line) of the distribution of  $\sqrt{n}(\hat{\theta} - \bar{\theta})$  for different sample sizes  $n$  (**Claw-Gaussian** model, 200 repetitions).

In the sequel, to lighten the notation, we assume without loss of generality that  $V = U$ . Thus, for all  $\theta \in U$ , we have  $\alpha(\theta) \in \Lambda^\circ$  and  $L(\theta, \alpha(\theta)) = \sup_{\alpha \in \Lambda} L(\theta, \alpha)$  (with  $\alpha(\bar{\theta}) = \bar{\alpha}$  by  $(H_1)$ ). Accordingly,  $\nabla_2 L(\theta, \alpha(\theta)) = 0, \forall \theta \in U$ . Also, since  $H_2 L(\bar{\theta}, \bar{\alpha})$  is invertible by  $(H_H)$  and since the function  $(\theta, \alpha) \mapsto H_2 L(\theta, \alpha)$  is continuous, there exists an open set  $U' \subset U$  such that  $H_2 L(\theta, \alpha)$  is invertible as soon as  $(\theta, \alpha) \in (U', \alpha(U'))$ . Without loss of generality, we assume that  $U' = U$ . Thus, by the chain rule, the function  $\alpha$  is of class  $C^2$  in a neighborhood  $U' \subset U$  of  $\bar{\theta}$ , say  $U' = U$ , with Jacobian matrix given by

$$J(\alpha)_\theta = -H_2 L(\theta, \alpha(\theta))^{-1} J(\nabla_2 L(\cdot, \alpha(\theta)))_\theta, \quad \forall \theta \in U.$$

We note that  $H_2 L(\theta, \alpha(\theta))^{-1}$  is of format  $q \times q$  and  $J(\nabla_2 L(\cdot, \alpha(\theta)))_\theta$  of format  $q \times p$ .

Now, for each  $\theta \in U$ , we let  $\hat{\alpha}(\theta)$  be such that  $\hat{L}(\theta, \hat{\alpha}(\theta)) = \sup_{\alpha \in \Lambda} \hat{L}(\theta, \alpha)$ .



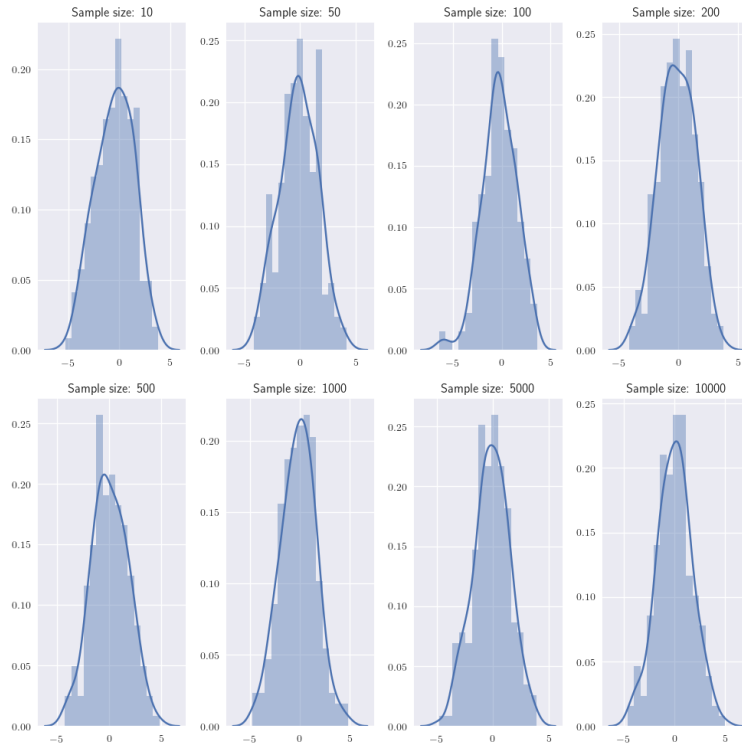


Fig 10: Histograms and kernel estimates (continuous line) of the distribution of  $\sqrt{n}(\hat{\theta} - \bar{\theta})$  for different sample sizes  $n$  (**Exponential-Uniform** model, 200 repetitions).

Clearly,

$$\begin{aligned}
& |L(\theta, \hat{\alpha}(\theta)) - L(\theta, \alpha(\theta))| \\
& \leq |L(\theta, \hat{\alpha}(\theta)) - \hat{L}(\theta, \hat{\alpha}(\theta))| + |\hat{L}(\theta, \hat{\alpha}(\theta)) - L(\theta, \alpha(\theta))| \\
& \leq \sup_{\alpha \in \Lambda} |L(\theta, \alpha) - \hat{L}(\theta, \alpha)| + \left| \sup_{\alpha \in \Lambda} \hat{L}(\theta, \alpha) - \sup_{\alpha \in \Lambda} L(\theta, \alpha) \right| \\
& \leq 2 \sup_{\alpha \in \Lambda} |\hat{L}(\theta, \alpha) - L(\theta, \alpha)|.
\end{aligned}$$

Therefore, by Lemma 4.1,  $\sup_{\theta \in U} |L(\theta, \hat{\alpha}(\theta)) - L(\theta, \alpha(\theta))| \rightarrow 0$  almost surely. The event on which this convergence holds does not depend upon  $\theta \in U$ , and, arguing as in the proof of Theorem 4.2, we deduce that under  $(H_1)$ ,  $\mathbb{P}(\hat{\alpha}(\theta) \rightarrow \alpha(\theta) \forall \theta \in U) = 1$ . Since  $\alpha(\theta) \in \Lambda^\circ$  for all  $\theta \in U$ , we also have  $\mathbb{P}(\hat{\alpha}(\theta) \in \Lambda^\circ \forall \theta \in U) \rightarrow 1$  as  $n \rightarrow \infty$ . Thus, in the sequel, it will be assumed without loss of generality that, for all  $\theta \in U$ ,  $\hat{\alpha}(\theta) \in \Lambda^\circ$ .

Still by Lemma 4.1,  $\sup_{\theta \in \Theta, \alpha \in \Lambda} \|H_2 \hat{L}(\theta, \alpha) - H_2 L(\theta, \alpha)\| \rightarrow 0$  almost surely. Since  $H_2 L(\theta, \alpha)$  is invertible on  $U \times \alpha(U)$ , we have

$$\mathbb{P}(H_2 \hat{L}(\theta, \alpha) \text{ invertible } \forall (\theta, \alpha) \in U \times \alpha(U)) \rightarrow 1.$$

Thus, we may and will assume that  $H_2 \hat{L}(\theta, \alpha)$  is invertible for all  $(\theta, \alpha) \in U \times \alpha(U)$ .

Next, since  $\hat{\alpha}(\theta) \in \Lambda^\circ$  for all  $\theta \in U$ , one has  $\nabla_2 \hat{L}(\theta, \hat{\alpha}(\theta)) = 0$ . Therefore, by the chain rule,  $\hat{\alpha}$  is of class  $C^2$  on  $U$ , with Jacobian matrix

$$J(\hat{\alpha})_\theta = -H_2 \hat{L}(\theta, \hat{\alpha}(\theta))^{-1} J(\nabla_2 \hat{L}(\cdot, \hat{\alpha}(\theta)))_\theta, \quad \forall \theta \in U.$$

Let  $\hat{V}(\theta) \stackrel{\text{def}}{=} \hat{L}(\theta, \hat{\alpha}(\theta)) = \sup_{\alpha \in \Lambda} \hat{L}(\theta, \alpha)$ . By the envelope theorem,  $\hat{V}$  is of class  $C^2$ ,  $\nabla \hat{V}(\theta) = \nabla_1 \hat{L}(\theta, \hat{\alpha}(\theta))$ , and

$$H \hat{V}(\theta) = H_1 \hat{L}(\theta, \hat{\alpha}(\theta)) + J(\nabla_1 \hat{L}(\theta, \cdot))_{\hat{\alpha}(\theta)} J(\hat{\alpha})_\theta.$$

Recall that  $\hat{\theta} \rightarrow \bar{\theta}$  almost surely by Theorem 4.2, so that we may assume that  $\hat{\theta} \in \Theta^\circ$  by  $(H_1)$ . Moreover, we can also assume that  $\hat{\theta} + t(\hat{\theta} - \bar{\theta}) \in U$ ,  $\forall t \in [0, 1]$ . Thus, by a Taylor series expansion with integral remainder, we have

$$(4.3) \quad 0 = \nabla \hat{V}(\hat{\theta}) = \nabla \hat{V}(\bar{\theta}) + \int_0^1 H \hat{V}(\bar{\theta} + t(\hat{\theta} - \bar{\theta})) dt (\hat{\theta} - \bar{\theta}).$$

Since  $\hat{\alpha}(\bar{\theta}) \in \Lambda^\circ$  and  $\hat{L}(\bar{\theta}, \hat{\alpha}(\bar{\theta})) = \sup_{\alpha \in \Lambda} \hat{L}(\bar{\theta}, \alpha)$ , one has  $\nabla_2 \hat{L}(\bar{\theta}, \hat{\alpha}(\bar{\theta})) = 0$ . Thus,

$$\begin{aligned} 0 &= \nabla_2 \hat{L}(\bar{\theta}, \hat{\alpha}(\bar{\theta})) \\ &= \nabla_2 \hat{L}(\bar{\theta}, \alpha(\bar{\theta})) + \int_0^1 H_2 \hat{L}(\bar{\theta}, \alpha(\bar{\theta}) + t(\hat{\alpha}(\bar{\theta}) - \alpha(\bar{\theta}))) dt (\hat{\alpha}(\bar{\theta}) - \alpha(\bar{\theta})). \end{aligned}$$

By Lemma 4.1, since  $\hat{\alpha}(\bar{\theta}) \rightarrow \alpha(\bar{\theta})$  almost surely, we have

$$\hat{I}_1 \stackrel{\text{def}}{=} \int_0^1 H_2 \hat{L}(\bar{\theta}, \alpha(\bar{\theta}) + t(\hat{\alpha}(\bar{\theta}) - \alpha(\bar{\theta}))) dt \rightarrow H_2 L(\bar{\theta}, \bar{\alpha}) \quad \text{almost surely.}$$

Because  $H_2 L(\bar{\theta}, \bar{\alpha})$  is invertible,  $\mathbb{P}(\hat{I}_1 \text{ invertible}) \rightarrow 1$  as  $n \rightarrow \infty$ . Therefore, we may assume, without loss of generality, that  $\hat{I}_1$  is invertible. Hence,

$$(4.4) \quad \hat{\alpha}(\bar{\theta}) - \alpha(\bar{\theta}) = -\hat{I}_1^{-1} \nabla_2 \hat{L}(\bar{\theta}, \alpha(\bar{\theta})).$$

Furthermore,

$$\nabla \hat{V}(\bar{\theta}) = \nabla_1 \hat{L}(\bar{\theta}, \hat{\alpha}(\bar{\theta})) = \nabla_1 \hat{L}(\bar{\theta}, \alpha(\bar{\theta})) + \hat{I}_2 (\hat{\alpha}(\bar{\theta}) - \alpha(\bar{\theta})),$$

where

$$\hat{I}_2 \stackrel{\text{def}}{=} \int_0^1 J(\nabla_1 \hat{L}(\bar{\theta}, \cdot))_{\alpha(\bar{\theta})+t(\hat{\alpha}(\bar{\theta})-\alpha(\bar{\theta}))} dt.$$

By Lemma 4.1,  $\hat{I}_2 \rightarrow J(\nabla_1 L(\bar{\theta}, \cdot))_{\alpha(\bar{\theta})}$  almost surely. Combining (4.3) and (4.4), we obtain

$$0 = \nabla_1 \hat{L}(\bar{\theta}, \alpha(\bar{\theta})) - \hat{I}_2 \hat{I}_1^{-1} \nabla_2 \hat{L}(\bar{\theta}, \alpha(\bar{\theta})) + \hat{I}_3 (\hat{\theta} - \bar{\theta}),$$

where

$$\hat{I}_3 \stackrel{\text{def}}{=} \int_0^1 H \hat{V}(\hat{\theta} + t(\hat{\theta} - \bar{\theta})) dt.$$

By technical Lemma 6.2, we have  $\hat{I}_3 \rightarrow HV(\bar{\theta})$  almost surely. So, by  $(H_V)$ , it can be assumed that  $\hat{I}_3$  is invertible. Consequently,

$$\hat{\theta} - \bar{\theta} = -\hat{I}_3^{-1} \nabla_1 \hat{L}(\bar{\theta}, \alpha(\bar{\theta})) + \hat{I}_3^{-1} \hat{I}_2 \hat{I}_1^{-1} \nabla_2 \hat{L}(\bar{\theta}, \alpha(\bar{\theta})),$$

or, equivalently, since  $\alpha(\bar{\theta}) = \bar{\alpha}$ ,

$$\hat{\theta} - \bar{\theta} = -\hat{I}_3^{-1} \nabla_1 \hat{L}(\bar{\theta}, \bar{\alpha}) + \hat{I}_3^{-1} \hat{I}_2 \hat{I}_1^{-1} \nabla_2 \hat{L}(\bar{\theta}, \bar{\alpha}).$$

Using Lemma 4.1, we conclude that  $\sqrt{n}(\hat{\theta} - \bar{\theta})$  has the same limit distribution as

$$S_n \stackrel{\text{def}}{=} -\sqrt{n}HV(\bar{\theta})^{-1} \nabla_1 \hat{L}(\bar{\theta}, \bar{\alpha}) + \sqrt{n}HV(\bar{\theta})^{-1} J(\nabla_1 L(\bar{\theta}, \cdot))_{\bar{\alpha}} H_2 L(\bar{\theta}, \bar{\alpha})^{-1} \nabla_2 \hat{L}(\bar{\theta}, \bar{\alpha}).$$

Let

$$\ell_i(\theta, \alpha) = \ln D_\alpha(X_i) + \ln(1 - D_\alpha \circ G_\theta(Z_i)), \quad 1 \leq i \leq n.$$

With this notation, we have

$$S_n = \frac{1}{\sqrt{n}} \sum_{i=1}^n \left( -HV(\bar{\theta})^{-1} \nabla_1 \ell_i(\bar{\theta}, \bar{\alpha}) + HV(\bar{\theta})^{-1} J(\nabla_1 L(\bar{\theta}, \cdot))_{\bar{\alpha}} H_2 L(\bar{\theta}, \bar{\alpha})^{-1} \nabla_2 \ell_i(\bar{\theta}, \bar{\alpha}) \right).$$

One has  $\nabla V(\bar{\theta}) = 0$ , since  $V(\bar{\theta}) = \inf_{\theta \in \Theta} V(\theta)$  and  $\bar{\theta} \in \Theta^\circ$ . Therefore, under  $(H'_{\text{reg}})$ ,  $\mathbb{E} \nabla_1 \ell_i(\bar{\theta}, \bar{\alpha}) = \nabla_1 \mathbb{E} \ell_i(\bar{\theta}, \bar{\alpha}) = \nabla_1 L(\bar{\theta}, \bar{\alpha}) = \nabla V(\bar{\theta}) = 0$ . Similarly, we have  $\mathbb{E} \nabla_2 \ell_i(\bar{\theta}, \bar{\alpha}) = \nabla_2 \mathbb{E} \ell_i(\bar{\theta}, \bar{\alpha}) = \nabla_2 L(\bar{\theta}, \bar{\alpha}) = 0$ , since  $L(\bar{\theta}, \bar{\alpha}) = \sup_{\alpha \in \Lambda} L(\bar{\theta}, \alpha)$  and  $\bar{\alpha} \in \Lambda^\circ$ . Using the central limit theorem, we conclude that

$$\sqrt{n}(\hat{\theta} - \bar{\theta}) \xrightarrow{\mathcal{L}} Z,$$

where  $Z$  is a Gaussian random variable with mean 0 and covariance matrix

$$\mathbf{V} = \text{Var} \left[ -HV(\bar{\theta})^{-1} \nabla_1 \ell_1(\bar{\theta}, \bar{\alpha}) + HV(\bar{\theta})^{-1} J(\nabla_1 L(\bar{\theta}, \cdot))_{\bar{\alpha}} H_2 L(\bar{\theta}, \bar{\alpha})^{-1} \nabla_2 \ell_1(\bar{\theta}, \bar{\alpha}) \right].$$

□

**5. Conclusion and perspectives.** In this paper, we have presented a theoretical study of the original Generative Adversarial Networks (GAN) algorithm, which consists in building a generative model of an unknown distribution from samples from that distribution. The key idea of the procedure is to simultaneously train the generative model (the generators) and an adversary (the discriminators) that tries to distinguish between real and generated samples. We made a small step towards a better understanding of this generative process by analyzing some optimality properties of the problem in terms of Jensen-Shannon divergence in Section 2, and explored the role of the discriminator family via approximation arguments in Section 3. Finally, taking a statistical view, we studied in Section 4 some large sample properties (convergence and asymptotic normality) of the parameter describing the empirically selected generator. Some numerical experiments were conducted to illustrate the results.

The point of view embraced in the article is statistical, in that it takes into account the variability of the data and its impact on the quality of the estimators. This point of view is different from the classical approach encountered in the literature on GANs, which mainly focuses on the effective computation of the parameters using optimization procedures. In this sense, our results must be thought of as a complementary insight. We realize however that the simplified context in which we have placed ourselves, as well as some of the assumptions we have made, are quite far from the typical situations in which GANs algorithms are used. Thus, our work should be seen as a first step towards a more realistic understanding of GANs, and certainly not as a definitive explanation for their excellent practical performance. We give below three avenues of theoretical research that we believe should be explored as a priority.

1. One of the basic assumptions is that the family of densities  $\{p_\theta\}_{\theta \in \Theta}$  (associated with the generators  $\{G_\theta\}_{\theta \in \Theta}$ ) and the unknown density  $p^*$  are dominated by the same measure  $\mu$  on the same subset  $E$  of  $\mathbb{R}^d$ . In a way, this means that we already have some kind of information on the support of  $p^*$ , which will typically be a manifold in  $\mathbb{R}^d$  of dimension smaller than  $d'$  (the dimension of  $Z$ ). Therefore, the random variable  $Z$ , the dimension  $d'$  of the so-called latent space  $\mathbb{R}^{d'}$ , and the parametric model  $\{G_\theta\}_{\theta \in \Theta}$  should be carefully tuned in order to match this constraint. From a practical perspective, the original article of [Goodfellow et al. \(2014\)](#) suggests using for  $Z$  a uniform or Gaussian distribution of small dimension, without further investigation. [Mirza and Osindero \(2014\)](#) and [Radford et al. \(2016\)](#), who have surprisingly good practical results with a deep convolutional generator, both use a 100-dimensional uniform distribution to represent respectively  $28 \times 28$  and  $64 \times 64$  pixel images. Many papers have been focusing on either decomposing the latent space  $\mathbb{R}^{d'}$  to force specified portions of this space to corre-

spond to different variations (as, e.g., in [Donahue et al., 2018](#)) or inverting the generators (e.g., [Lipton and Tripathi, 2017](#); [Srivastava et al., 2017](#); [Bojanowski et al., 2018](#)). However, to the best of our knowledge, there is to date no theoretical result tackling the impact of  $d'$  and  $Z$  on the performance of GANs, and it is our belief that a thorough mathematical investigation of this issue is needed for a better understanding of the generating process. Similarly, whenever the  $\{G_\theta\}_{\theta \in \Theta}$  are neural networks, the link between the networks (number of layers, dimensionality of  $\Theta$ , etc.) and the target  $p^*$  (support, dominating measure, etc.) is also a fundamental question, which should be addressed at a theoretical level.

2. Assumptions  $(H_\varepsilon)$  and  $(H'_\varepsilon)$  highlight the essential role played by the discriminators to approximate the optimal functions  $D_\theta^*$ . We believe that this point is critical for the theoretical analysis of GANs, and that it should be further developed in the context of neural networks, with a potentially large number of hidden layers.
3. Theorem 4.2 (convergence of the estimated parameter) and Theorem 4.3 (asymptotic normality) hold under the assumption that the model is identifiable (uniqueness of  $\bar{\theta}$  and  $\bar{\alpha}$ ). This identifiability assumption is hardly satisfied in the high-dimensional context of (deep) neural networks, where the function to be optimized displays a very wild landscape, without immediate convexity or concavity. Thus, to take one more step towards a more realistic model, it would be interesting to shift the parametric point of view and move towards results concerning the convergence of distributions not parameters.

## 6. Technical results.

6.1. *Proof of Theorem 3.1.* Let  $\varepsilon \in (0, 1/(2M))$ ,  $m \in (0, 1/2)$ , and  $D \in \mathcal{D}$  be such that  $m \leq D \leq 1 - m$  and  $\|D - D_{\bar{\theta}}^*\|_2 \leq \varepsilon$ . Observe that

$$\begin{aligned}
 L(\bar{\theta}, D) &= \int \ln(D) p^* d\mu + \int \ln(1-D) p_{\bar{\theta}} d\mu \\
 (6.1) \quad &= \int \ln\left(\frac{D}{D_{\bar{\theta}}^*}\right) p^* d\mu + \int \ln\left(\frac{1-D}{1-D_{\bar{\theta}}^*}\right) p_{\bar{\theta}} d\mu + 2D_{\text{JS}}(p^*, p_{\bar{\theta}}) - \ln 4.
 \end{aligned}$$

We first derive a lower bound on the quantity

$$\begin{aligned}
 I &\stackrel{\text{def}}{=} \int \ln\left(\frac{D}{D_{\bar{\theta}}^*}\right) p^* d\mu + \int \ln\left(\frac{1-D}{1-D_{\bar{\theta}}^*}\right) p_{\bar{\theta}} d\mu \\
 &= \int \ln\left(\frac{D(p^* + p_{\bar{\theta}})}{p^*}\right) p^* d\mu + \int \ln\left(\frac{(1-D)(p^* + p_{\bar{\theta}})}{p_{\bar{\theta}}}\right) p_{\bar{\theta}} d\mu.
 \end{aligned}$$

Let  $dP^* = p^*d\mu$ ,  $dP_{\bar{\theta}} = p_{\bar{\theta}}d\mu$ ,

$$d\kappa = \frac{D(p^* + p_{\bar{\theta}})}{\int D(p^* + p_{\bar{\theta}})d\mu}d\mu, \quad \text{and} \quad d\kappa' = \frac{(1-D)(p^* + p_{\bar{\theta}})}{\int (1-D)(p^* + p_{\bar{\theta}})d\mu}d\mu.$$

Observe, since  $m \leq D \leq 1 - m$ , that  $P^* \ll \kappa$  and  $P_{\bar{\theta}} \ll \kappa'$ . With this notation, we have

$$(6.2) \quad I = -D_{\text{KL}}(P^* \parallel \kappa) - D_{\text{KL}}(P_{\bar{\theta}} \parallel \kappa') + \ln \left[ \int D(p^* + p_{\bar{\theta}})d\mu (2 - \int D(p^* + p_{\bar{\theta}})d\mu) \right].$$

Since

$$\int D(p^* + p_{\bar{\theta}})d\mu = \int (D - D_{\bar{\theta}}^*)(p^* + p_{\bar{\theta}})d\mu + 1,$$

the Cauchy-Schwartz inequality leads to

$$(6.3) \quad \left| \int D(p^* + p_{\bar{\theta}})d\mu - 1 \right| \leq \|D - D_{\bar{\theta}}^*\|_2 \|p^* + p_{\bar{\theta}}\|_2 \leq 2M\varepsilon,$$

because both  $p^*$  and  $p_{\bar{\theta}}$  are bounded by  $M$ . Thus,

$$(6.4) \quad \ln \left[ \int D(p^* + p_{\bar{\theta}})d\mu (2 - \int D(p^* + p_{\bar{\theta}})d\mu) \right] \geq \ln(1 - 4M^2\varepsilon^2) \geq -\frac{4M^2\varepsilon^2}{1 - 4M^2\varepsilon^2},$$

using the inequality  $\ln(1 - x) \geq -x/(1 - x)$  for  $x \in [0, 1)$ . Moreover, recalling that the Kullback-Leibler divergence is smaller than the chi-square divergence, and letting  $\bar{F} = F/(\int Fd\mu)$  for  $F \in L^1(\mu)$ , we have

$$D_{\text{KL}}(P^* \parallel \kappa) \leq \int \left( \frac{p^*}{D(p^* + p_{\bar{\theta}})} - 1 \right)^2 \frac{1}{D(p^* + p_{\bar{\theta}})}d\mu.$$

Hence, letting  $J \stackrel{\text{def}}{=} \int D(p^* + p_{\bar{\theta}})d\mu$ , we see that

$$\begin{aligned} & D_{\text{KL}}(P^* \parallel \kappa) \\ & \leq \frac{1}{J} \int \left( p^* \int D(p^* + p_{\bar{\theta}})d\mu - D(p^* + p_{\bar{\theta}}) \right)^2 \frac{1}{D(p^* + p_{\bar{\theta}})}d\mu \\ & = \frac{1}{J} \int \left( p^* \int (D - D_{\bar{\theta}}^*)(p^* + p_{\bar{\theta}})d\mu + (D_{\bar{\theta}}^* - D)(p^* + p_{\bar{\theta}}) \right)^2 \frac{1}{D(p^* + p_{\bar{\theta}})}d\mu. \end{aligned}$$

Since  $\varepsilon < 1/(2M)$ , inequality (6.3) gives  $1/J \leq c_1$  for some constant  $c_1 > 0$ . By Cauchy-Schwarz and  $(a+b)^2 \leq 2(a^2+b^2)$ , we obtain

$$\begin{aligned} D_{\text{KL}}(P^* \parallel \kappa) &\leq 2c_1 \left( \int \left( \int (D - D_{\bar{\theta}}^*)(p^* + p_{\bar{\theta}}) d\mu \right)^2 \frac{(p^*)^2}{D(p^* + p_{\bar{\theta}})} d\mu + \int (D_{\bar{\theta}}^* - D)^2 \frac{p^* + p_{\bar{\theta}}}{D} d\mu \right) \\ &\leq 2c_1 \left( \|D - D_{\bar{\theta}}^*\|_2^2 \|p^* + p_{\bar{\theta}}\|_2^2 \int \frac{(p^*)^2}{D(p^* + p_{\bar{\theta}})} d\mu + \int (D_{\bar{\theta}}^* - D)^2 \frac{p^* + p_{\bar{\theta}}}{D} d\mu \right). \end{aligned}$$

Therefore, since  $p^* \leq M$ ,  $p_{\bar{\theta}} \leq M$ , and  $D \geq m$ ,

$$D_{\text{KL}}(P^* \parallel \kappa) \leq 2c_1 \left( \frac{4M^2}{m} + \frac{2M}{m} \right) \varepsilon^2.$$

One proves with similar arguments that

$$D_{\text{KL}}(P_{\bar{\theta}} \parallel \kappa') \leq 2c_1 \left( \frac{4M^2}{m} + \frac{2M}{m} \right) \varepsilon^2.$$

Combining these two inequalities with (6.2) and (6.4), we see that  $I \geq -c_2 \varepsilon^2$  for some constant  $c_2 > 0$  that depends only upon  $M$  and  $m$ . Getting back to identity (6.1), we conclude that

$$2D_{\text{JS}}(p^*, p_{\bar{\theta}}) \leq c_2 \varepsilon^2 + L(\bar{\theta}, D) + \ln 4.$$

But

$$\begin{aligned} L(\bar{\theta}, D) &\leq \sup_{D \in \mathcal{D}} L(\bar{\theta}, D) \leq \sup_{D \in \mathcal{D}} L(\theta^*, D) \\ &\quad \text{(by definition of } \bar{\theta} \text{)} \\ &\leq \sup_{D \in \mathcal{D}_\infty} L(\theta^*, D) \\ &= L(\theta^*, D_{\theta^*}^*) = 2D_{\text{JS}}(p^*, p_{\theta^*}) - \ln 4. \end{aligned}$$

Thus,

$$2D_{\text{JS}}(p^*, p_{\bar{\theta}}) \leq c_2 \varepsilon^2 + 2D_{\text{JS}}(p^*, p_{\theta^*}).$$

This shows the right-hand side of inequality (3.1). To prove the left-hand side, just note that by inequality (2.2),

$$D_{\text{JS}}(p^*, p_{\theta^*}) \leq D_{\text{JS}}(p^*, p_{\bar{\theta}}).$$

6.2. *Proof of Lemma 4.1.* To simplify the notation, we set

$$\Delta = \frac{\partial^{a+b+c+d}}{\partial \theta_i^a \partial \theta_j^b \partial \alpha_\ell^c \partial \alpha_m^d}.$$

Using McDiarmid's inequality (McDiarmid, 1989), we see that there exists a constant  $c > 0$  such that, for all  $\varepsilon > 0$ ,

$$\mathbb{P}\left(\left|\sup_{\theta \in \Theta, \alpha \in \Lambda} |\Delta \hat{L}(\theta, \alpha) - \Delta L(\theta, \alpha)| - \mathbb{E} \sup_{\theta \in \Theta, \alpha \in \Lambda} |\Delta \hat{L}(\theta, \alpha) - \Delta L(\theta, \alpha)|\right| \geq \varepsilon\right) \leq 2e^{-c\varepsilon^2}.$$

Therefore, by the Borel-Cantelli lemma,

$$(6.5) \quad \sup_{\theta \in \Theta, \alpha \in \Lambda} |\Delta \hat{L}(\theta, \alpha) - \Delta L(\theta, \alpha)| - \mathbb{E} \sup_{\theta \in \Theta, \alpha \in \Lambda} |\Delta \hat{L}(\theta, \alpha) - \Delta L(\theta, \alpha)| \rightarrow 0$$

almost surely. It is also easy to verify that under Assumptions  $(H'_{\text{reg}})$ , the process  $(\Delta \hat{L}(\theta, \alpha) - \Delta L(\theta, \alpha))_{\theta \in \Theta, \alpha \in \Lambda}$  is subgaussian. Thus, as in the proof of Theorem 4.1, we obtain via Dudley's inequality that

$$(6.6) \quad \mathbb{E} \sup_{\theta \in \Theta, \alpha \in \Lambda} |\Delta \hat{L}(\theta, \alpha) - \Delta L(\theta, \alpha)| = O\left(\frac{1}{\sqrt{n}}\right),$$

since  $\mathbb{E} \Delta \hat{L}(\theta, \alpha) = \Delta L(\theta, \alpha)$ . The result follows by combining (6.5) and (6.6).

### 6.3. Some technical lemmas.

LEMMA 6.1. *Under Assumptions  $(H'_{\text{reg}})$  and  $(H_1)$ , there exists an open set  $V \subset \Theta^\circ$  containing  $\bar{\theta}$  such that, for all  $\theta \in V$ ,  $\arg \max_{\alpha \in \Lambda} L(\theta, \alpha) \cap \Lambda^\circ \neq \emptyset$ .*

PROOF. Assume that the statement is not true. Then there exists a sequence  $(\theta_k)_k \subset \Theta$  such that  $\theta_k \rightarrow \bar{\theta}$  and, for all  $k$ ,  $\alpha_k \in \partial \Lambda$ , where  $\alpha_k \in \arg \max_{\alpha \in \Lambda} L(\theta_k, \alpha)$ . Thus, since  $\Lambda$  is compact, even if this means extracting a subsequence, one has  $\alpha_k \rightarrow z \in \partial \Lambda$  as  $k \rightarrow \infty$ . By the continuity of  $L$ ,  $L(\bar{\theta}, \alpha_k) \rightarrow L(\bar{\theta}, z)$ . But

$$\begin{aligned} |L(\bar{\theta}, \alpha_k) - L(\bar{\theta}, \bar{\alpha})| &\leq |L(\bar{\theta}, \alpha_k) - L(\theta_k, \alpha_k)| + |L(\theta_k, \alpha_k) - L(\bar{\theta}, \bar{\alpha})| \\ &\leq \sup_{\alpha \in \Lambda} |L(\bar{\theta}, \alpha) - L(\theta_k, \alpha)| + \left| \sup_{\alpha \in \Lambda} L(\theta_k, \alpha) - \sup_{\alpha \in \Lambda} L(\bar{\theta}, \alpha) \right| \\ &\leq 2 \sup_{\alpha \in \Lambda} |L(\bar{\theta}, \alpha) - L(\theta_k, \alpha)|, \end{aligned}$$

which tends to zero as  $k \rightarrow \infty$  by  $(H'_D)$  and  $(H'_p)$ . Therefore,  $L(\bar{\theta}, z) = L(\bar{\theta}, \bar{\alpha})$  and, in turn,  $z = \bar{\alpha}$  by  $(H_1)$ . Since  $z \in \partial \Lambda$  and  $\bar{\alpha} \in \Lambda^\circ$ , this is a contradiction.  $\square$

LEMMA 6.2. *Under Assumptions  $(H'_{\text{reg}})$ ,  $(H_1)$ , and  $(H_{\text{loc}})$ , one has  $\hat{I}_3 \rightarrow HV(\bar{\theta})$  almost surely.*



PROOF. We have

$$\hat{I}_3 = \int_0^1 H\hat{V}(\hat{\theta} + t(\hat{\theta} - \bar{\theta}))dt = \int_0^1 (H_1\hat{L}(\hat{\theta}_t, \hat{\alpha}(\hat{\theta}_t)) + J(\nabla_1\hat{L}(\hat{\theta}_t, \cdot))_{\hat{\alpha}(\hat{\theta}_t)}J(\hat{\alpha})_{\hat{\theta}_t})dt,$$

where we set  $\hat{\theta}_t = \hat{\theta} + t(\hat{\theta} - \bar{\theta})$ . Note that  $\hat{\theta}_t \in U$  for all  $t \in [0, 1]$ . By Lemma 4.1,

$$\begin{aligned} & \sup_{t \in [0, 1]} \|H_1\hat{L}(\hat{\theta}_t, \hat{\alpha}(\hat{\theta}_t)) - H_1L(\hat{\theta}_t, \hat{\alpha}(\hat{\theta}_t))\| \\ & \leq \sup_{\theta \in \Theta, \alpha \in \Lambda} \|H_1\hat{L}(\theta, \alpha) - H_1L(\theta, \alpha)\| \rightarrow 0 \quad \text{almost surely.} \end{aligned}$$

Also, by Theorem 4.2, for all  $t \in [0, 1]$ ,  $\hat{\theta}_t \rightarrow \bar{\theta}$  almost surely. Besides,

$$\begin{aligned} |L(\bar{\theta}, \hat{\alpha}(\hat{\theta}_t)) - L(\bar{\theta}, \alpha(\bar{\theta}))| & \leq |L(\bar{\theta}, \hat{\alpha}(\hat{\theta}_t)) - L(\hat{\theta}_t, \hat{\alpha}(\hat{\theta}_t))| \\ & \quad + |L(\hat{\theta}_t, \hat{\alpha}(\hat{\theta}_t)) - L(\bar{\theta}, \alpha(\bar{\theta}))| \\ & \leq \sup_{\alpha \in \Lambda} |L(\bar{\theta}, \alpha) - L(\hat{\theta}_t, \alpha)| \\ & \quad + 2 \sup_{\theta \in \Theta, \alpha \in \Lambda} |\hat{L}(\theta, \alpha) - L(\theta, \alpha)|. \end{aligned}$$

Thus, via  $(H'_{\text{reg}})$ ,  $(H_1)$ , and Lemma 4.1, we conclude that almost surely, for all  $t \in [0, 1]$ , one has  $\hat{\alpha}(\hat{\theta}_t) \rightarrow \alpha(\bar{\theta}) = \bar{\alpha}$ . Accordingly, almost surely, for all  $t \in [0, 1]$ ,  $H_1L(\hat{\theta}_t, \hat{\alpha}(\hat{\theta}_t)) \rightarrow H_1L(\bar{\theta}, \bar{\alpha})$ . Since  $H_1L(\theta, \alpha)$  is bounded under  $(H'_D)$  and  $(H'_p)$ , the Lebesgue dominated convergence theorem leads to

$$(6.7) \quad \int_0^1 H_1\hat{L}(\hat{\theta}_t, \hat{\alpha}(\hat{\theta}_t))dt \rightarrow H_1L(\bar{\theta}, \bar{\alpha}) \quad \text{almost surely.}$$

Furthermore,

$$J(\hat{\alpha})_{\theta} = -H_2\hat{L}(\theta, \hat{\alpha}(\theta))^{-1}J(\nabla_2\hat{L}(\cdot, \hat{\alpha}(\theta)))_{\theta}, \quad \forall(\theta, \alpha) \in U \times \alpha(U),$$

where  $U$  is the open set defined in the proof of Theorem 4.3. By the cofactor method,  $H_2\hat{L}(\theta, \alpha)^{-1}$  takes the form

$$H_2\hat{L}(\theta, \alpha)^{-1} = \frac{\hat{c}(\theta, \alpha)}{\det(H_2\hat{L}(\theta, \alpha))},$$

where  $\hat{c}(\theta, \alpha)$  is the matrix of cofactors associated with  $H_2\hat{L}(\theta, \alpha)$ . Thus, each component of  $-H_2\hat{L}(\theta, \alpha)^{-1}J(\nabla_2\hat{L}(\cdot, \alpha))_{\theta}$  is a quotient of a multilinear form of the partial derivatives of  $\hat{L}$  evaluated at  $(\theta, \alpha)$  divided by  $\det(H_2\hat{L}(\theta, \alpha))$ , which is itself a multilinear form in the  $\frac{\partial^2 \hat{L}}{\partial \alpha_i \partial \alpha_j}(\theta, \alpha)$ . Hence, by Lemma 4.1, we have

$$\sup_{\theta \in U, \alpha \in \alpha(U)} \|H_2\hat{L}(\theta, \alpha)^{-1}J(\nabla_2\hat{L}(\cdot, \alpha))_{\theta} - H_2L(\theta, \alpha)^{-1}J(\nabla_2L(\cdot, \alpha))_{\theta}\| \rightarrow 0$$

almost surely. So, for all  $n$  large enough,

$$\begin{aligned} & \sup_{t \in [0,1]} \|J(\hat{\alpha})_{\hat{\theta}_t} + H_2L(\hat{\theta}_t, \hat{\alpha}(\hat{\theta}_t))^{-1}J(\nabla_2L(\cdot, \hat{\alpha}(\hat{\theta}_t)))_{\hat{\theta}_t}\| \\ & \leq \sup_{\theta \in U, \alpha \in \alpha(U)} \|H_2\hat{L}(\theta, \alpha)^{-1}J(\nabla_2\hat{L}(\cdot, \alpha))_{\theta} - H_2L(\theta, \alpha)^{-1}J(\nabla_2L(\cdot, \alpha))_{\theta}\| \\ & \rightarrow 0 \quad \text{almost surely.} \end{aligned}$$

We know that almost surely, for all  $t \in [0, 1]$ ,  $\hat{\alpha}(\hat{\theta}_t) \rightarrow \bar{\alpha}$ . Thus, since the function  $U \times \alpha(U) \ni (\theta, \alpha) \mapsto H_2L(\theta, \alpha)^{-1}J(\nabla_2L(\cdot, \alpha))_{\theta}$  is continuous, we have almost surely, for all  $t \in [0, 1]$ ,

$$H_2\hat{L}(\hat{\theta}_t, \hat{\alpha}(\hat{\theta}_t))^{-1}J(\nabla_2\hat{L}(\cdot, \hat{\alpha}(\hat{\theta}_t)))_{\hat{\theta}_t} \rightarrow H_2L(\bar{\theta}, \bar{\alpha})^{-1}J(\nabla_2L(\cdot, \bar{\alpha}))_{\bar{\theta}}.$$

Therefore, almost surely, for all  $t \in [0, 1]$ ,  $J(\hat{\alpha})_{\hat{\theta}_t} \rightarrow J(\alpha)_{\bar{\theta}}$ . Similarly, almost surely, for all  $t \in [0, 1]$ ,  $J(\nabla_1\hat{L}(\hat{\theta}_t, \cdot))_{\hat{\alpha}(\hat{\theta}_t)} \rightarrow J(\nabla_1L(\bar{\theta}, \cdot))_{\bar{\alpha}}$ . All involved quantities are uniformly bounded in  $t$ , and so, by the Lebesgue dominated convergence theorem, we conclude that

$$(6.8) \quad \int_0^1 J(\nabla_1\hat{L}(\hat{\theta}_t, \cdot))_{\hat{\alpha}(\hat{\theta}_t)} J(\hat{\alpha})_{\hat{\theta}_t} dt \rightarrow J(\nabla_1L(\bar{\theta}, \cdot))_{\bar{\alpha}} J(\alpha)_{\bar{\theta}} \quad \text{almost surely.}$$

Consequently, by combining (6.7) and (6.8),

$$\hat{I}_3 \rightarrow H_1L(\bar{\theta}, \bar{\alpha}) + J(\nabla_1L(\bar{\theta}, \cdot))_{\bar{\alpha}} J(\alpha)_{\bar{\theta}} = HV(\bar{\theta}) \quad \text{almost surely,}$$

as desired.  $\square$

## ACKNOWLEDGEMENTS

We thank Flavian Vasile (Criteo AI Lab) and Antoine Picard-Weibel (ENS Ulm) for stimulating discussions and insightful suggestions. We also thank the Associate Editor and two anonymous referees for their careful reading of the paper and constructive comments, which led to a substantial improvement of the document.

## REFERENCES

- M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G.S. Corrado, A. Davis, J. Dean, M. Devin, S. Ghemawat, I. Goodfellow, A. Harp, G. Irving, M. Isard, Y. Jia, R. Jozefowicz, L. Kaiser, M. Kudlur, J. Levenberg, D. Mané, R. Monga, S. Moore, D. Murray, C. Olah, M. Schuster, J. Shlens, B. Steiner, I. Sutskever, K. Talwar, P. Tucker, V. Vanhoucke, V. Vasudevan, F. Viégas, O. Vinyals, P. Warden, M. Wattenberg, M. Wicke, Y. Yu, and X. Zheng. TensorFlow: A system for large-scale machine learning, 2015. URL <https://www.tensorflow.org/>. Software available from tensorflow.org.
- T. Angles and S. Mallat. Generative networks as inverse problems with scattering transforms. In *International Conference on Learning Representations*, 2018.

- M. Arjovsky and L. Bottou. Towards principled methods for training generative adversarial networks. In *International Conference on Learning Representations*, 2017.
- M. Arjovsky, S. Chintala, and L. Bottou. Wasserstein generative adversarial networks. In D. Precup and Y. Whye Teh, editors, *Proceedings of the 34th International Conference on Machine Learning*, volume 70, pages 214–223. Proceedings of Machine Learning Research, 2017.
- P. Bojanowski, A. Joulin, D. Lopez-Paz, and A. Szlam. Optimizing the latent space of generative networks. In J. Dy and A. Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80, pages 600–609. Proceedings of Machine Learning Research, 2018.
- L. Devroye. Universal smoothing factor selection in density estimation: Theory and practice. *TEST*, 6:223–320, 1997.
- C. Donahue, Z.C. Lipton, A. Balsubramani, and J. McAuley. Semantically decomposing the latent spaces of generative adversarial networks. In *International Conference on Learning Representations*, 2018.
- G.K. Dziugaite, D.M. Roy, and Z. Ghahramani. Training generative neural networks via maximum mean discrepancy optimization. In M. Meila and T. Heskes, editors, *Proceedings of the Thirty-First Conference on Uncertainty in Artificial Intelligence*, pages 258–267. AUAI Press, Arlington, 2015.
- D.M. Endres and J.E. Schindelin. A new metric for probability distributions. *IEEE Transactions on Information Theory*, 49:1858–1860, 2003.
- I. Goodfellow. *NIPS 2016 Tutorial: Generative Adversarial Networks*. arXiv:1701.00160, 2016.
- I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and J. Bengio. Generative adversarial nets. In Z. Ghahramani, M. Welling, C. Cortes, N.D. Lawrence, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems 27*, pages 2672–2680. Curran Associates, Inc., Red Hook, 2014.
- T. Liang. *On how well generative adversarial networks learn densities: Nonparametric and parametric results*. arXiv:1811.03179, 2018.
- Z.C. Lipton and S. Tripathi. *Precise recovery of latent vectors from generative adversarial networks*. arxiv:1702.04782, 2017.
- S. Liu, O. Bousquet, and K. Chaudhuri. Approximation and convergence properties of generative adversarial learning. In I. Guyon, U.V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 5551–5559. Curran Associates, Inc., Red Hook, 2017.
- C. McDiarmid. On the method of bounded differences. In J. Siemons, editor, *Surveys in Combinatorics*, London Mathematical Society Lecture Note Series 141, pages 148–188. Cambridge University Press, Cambridge, 1989.
- M. Mirza and S. Osindero. *Conditional generative adversarial nets*. arXiv:1411.1784, 2014.
- S. Nowozin, B. Cseke, and R. Tomioka. f-GAN: Training generative neural samplers using variational divergence minimization. In D.D. Lee, M. Sugiyama, U.V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems 29*, pages 271–279. Curran Associates, Inc., Red Hook, 2016.
- A. Radford, L. Metz, and S. Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. In *International Conference on Learning Representations*, 2016.
- T. Salimans, I. Goodfellow, W. Zaremba, V. Cheung, A. Radford, and X. Chen. Improved techniques for training GANs. In D.D. Lee, M. Sugiyama, U.V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems 29*, pages 2234–2242. Curran Associates, Inc., Red Hook, 2016.
- A. Srivastava, L. Valkoz, C. Russell, M.U. Gutmann, and C. Sutton. Veegan: Reducing mode collapse in GANs using implicit variational learning. In I. Guyon, U.V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing*

*Systems 30*, pages 3308–3318. Curran Associates, Inc., Red Hook, 2017.

R. van Handel. *Probability in High Dimension*. APC 550 Lecture Notes, Princeton University, 2016.

P. Zhang, Q. Liu, D. Zhou, T. Xu, and X. He. On the discriminative-generalization tradeoff in GANs.

In *International Conference on Learning Representations*, 2018.

G. BIAU  
M. SANGNIER  
SORBONNE UNIVERSITÉ  
LABORATOIRE DE PROBABILITÉS,  
STATISTIQUE ET MODÉLISATION  
BOÎTE 158, TOUR 15-25  
4 PLACE JUSSIEU  
75005 PARIS, FRANCE  
E-MAIL: [gerard.biau@sorbonne-universite.fr](mailto:gerard.biau@sorbonne-universite.fr)  
[maxime.sangnier@sorbonne-universite.fr](mailto:maxime.sangnier@sorbonne-universite.fr)

B. CADRE  
UNIV RENNES, IRMAR  
ENS RENNES  
AVENUE ROBERT SCHUMANN  
35170 BRUZ, FRANCE  
E-MAIL: [benoit.cadre@ens-rennes.fr](mailto:benoit.cadre@ens-rennes.fr)

U. TANIELIAN  
CRITEO AI LAB  
32 RUE BLANCHE  
75009 PARIS, FRANCE  
E-MAIL: [u.tanielian@criteo.com](mailto:u.tanielian@criteo.com)