



On the robustness of randomized classifiers to adversarial examples

Rafael Pinot¹ · Laurent Meunier^{2,3} · Florian Yger³ · Cédric Gouy-Pailler⁴ · Yann Chevaleyre³ · Jamal Atif³

Received: 26 October 2021 / Revised: 9 June 2022 / Accepted: 13 June 2022 /
Published online: 2 August 2022
© The Author(s) 2022

Abstract

This paper investigates the theory of robustness against adversarial attacks. We focus on randomized classifiers (i.e. classifiers that output random variables) and provide a thorough analysis of their behavior through the lens of statistical learning theory and information theory. To this aim, we introduce a new notion of robustness for randomized classifiers, enforcing local Lipschitzness using probability metrics. Equipped with this definition, we make two new contributions. The first one consists in devising a new upper bound on the adversarial generalization gap of randomized classifiers. More precisely, we devise bounds on the generalization gap and the adversarial gap (i.e. the gap between the risk and the worst-case risk under attack) of randomized classifiers. The second contribution presents a yet simple but efficient noise injection method to design robust randomized classifiers. We show that our results are applicable to a wide range of machine learning models under mild hypotheses. We further corroborate our findings with experimental results using deep neural networks on standard image datasets, namely CIFAR-10 and CIFAR-100. On these tasks, we manage to design robust models that simultaneously achieve state-of-the-art accuracy (over 0.82 clean accuracy on CIFAR-10) and enjoy *guaranteed* robust accuracy bounds (0.45 against ℓ_2 adversaries with magnitude 0.5 on CIFAR-10).

Keywords Adversarial examples · Randomized classifier · Adversarial generalization gap · Information theory · Randomized smoothing

Editor: Willem Waegeman.

Rafael Pinot and Laurent Meunier have contributed equally to this work.

✉ Rafael Pinot
rafael.pinot@epfl.ch

✉ Laurent Meunier
laurent.meunier@dauphine.eu

¹ EPFL, 1015 Lausanne, Switzerland

² Meta AI Research, 75002 Paris, France

³ LAMSADE, CNRS, Université Paris-Dauphine, PSL Research University, 75016 Paris, France

⁴ CEA, List, Université Paris-Saclay, 91120 Palaiseau, France

1 Introduction

In the last few years, there has been a growing concern on adversarial example attacks in machine learning. An adversarial attack refers to a small (humanly imperceptible) change of an input specifically designed to fool a machine learning model. These attacks have recently come to light thanks to works by Biggio et al. (2013) and Szegedy et al. (2014) studying deep neural networks for image classification, although it was an existing topic in spam filter analysis (Dalvi et al., 2004; Lowd & Meeck, 2005; Globerson & Roweis, 2006). The vulnerability of state-of-the-art classifiers to these attacks has genuine security implications especially for deep neural networks used in AI-driven technologies such as self-driving cars, as repetitively demonstrated by Sharif et al. (2016), Sitawarin et al. (2018) and Yao et al. (2020). Besides security issues, this shows how little we know about the worst-case behaviors of models the industry uses daily. It is essential for the community to understand the very nature of this phenomenon in order to mitigate the threat.

Accordingly, a large body of works has been trying to design new models that would be less vulnerable to the adversarial setting (Goodfellow et al., 2015; Metzen et al., 2017; Xie et al., 2018; Hu et al., 2019; Verma & Swami, 2019) but most of them were proven (in time) to offer only limited protection against more sophisticated attacks (Carlini & Wagner, 2017; He et al., 2017; Athalye et al., 2018; Croce & Hein, 2020; Tramer et al., 2020). Among the defense strategies, randomization has proven effective in some contexts (Xie et al., 2018; Dhillon et al., 2018; Liu et al., 2018; He et al., 2019). Albeit these significant efforts, randomization techniques lack theoretical arguments. In this paper, we generalize the prior results from Pinot et al. (2019) by studying a general class of randomized classifiers, including randomized neural networks, for which we demonstrate adversarial robustness guarantees and analyze their generalization properties (see Sect. 2.3 for more details).

1.1 Supervised learning for image classification

Let us consider the supervised classification problem with an input space \mathcal{X} and an output space \mathcal{Y} . In the following, w.l.o.g. we will consider $\mathcal{X} \subset [-1, 1]^d$ to be a set of images, and $\mathcal{Y} := [K] := \{1, \dots, K\}$ a set of labels describing them. The goal of a supervised machine learning algorithm is to design classifier that maps any image $\mathbf{x} \in \mathcal{X}$ to a label $y \in \mathcal{Y}$. To do so, the learner has access to a *training sample* of n image-label pairs $\mathcal{S} := \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$. Each training pair (\mathbf{x}_i, y_i) is assumed to be drawn i.i.d. from a ground-truth distribution \mathcal{D} . To build a classifier, the usual strategy is to select a hypothesis function $\mathbf{h} : \mathcal{X} \rightarrow \mathcal{Y}$ from a pre-defined hypothesis class \mathcal{H} to minimize the *risk* with respect to \mathcal{D} . This risk minimization problem writes

$$\inf_{\mathbf{h} \in \mathcal{H}} \mathcal{R}(\mathbf{h}) := \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [\mathcal{L}_{0/1}(\mathbf{h}(\mathbf{x}), y)], \quad (1)$$

where $\mathcal{L}_{0/1}$, the 0/1 loss, outputs 1 when $\mathbf{h}(\mathbf{x}) \neq y$, and zero otherwise.

In practice, the learner does not have access to the ground-truth distribution; hence it cannot estimate the risk $\mathcal{R}(\mathbf{h})$. To find an approximate solution for Problem (1), a learning algorithm solves the *empirical risk minimization* problem instead. In this case, we simply replace the risk by its empirical counterpart over the training sample $\mathcal{S} := \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$. The empirical risk minimization problem writes

$$\inf_{\mathbf{h} \in \mathcal{H}} \mathcal{R}_S(\mathbf{h}) := \frac{1}{n} \sum_{i=1}^n \mathcal{L}_{0/1}(\mathbf{h}(\mathbf{x}_i), y_i). \quad (2)$$

Then, to evaluate how far the selected hypothesis is from the optimum, one wants to upper bound the difference between the risk and the empirical risk of any $\mathbf{h} \in \mathcal{H}$. This difference is known as the *generalization gap*.

1.2 Classification in the presence of an adversary

Given a hypothesis $\mathbf{h} \in \mathcal{H}$ and a sample $(x, y) \sim \mathcal{D}$, the goal of an adversary is to find a perturbation $\tau \in \mathcal{X}$ such that the following assertions *both* hold. First, the perturbation is imperceptible to humans. This means that a human cannot visually distinguish the standard example x from the *adversarial example* $x + \tau$. Second, the perturbation modifies x enough to make the classifier misclassify. More formally, the adversary seeks a perturbation $\tau \in \mathcal{X}$ such that $\mathbf{h}(x + \tau) \neq y$.

Although the notion of imperceptible modification is very natural for humans, it is genuinely hard to formalize. Despite these difficulties, in the image classification setting, a sufficient condition to ensure that the attack will remain undetected is to constrain the perturbation τ to have a small ℓ_p norm. This means that for any $p \in [1, \infty]$, there exists a threshold $\alpha_p > 0$ for which any perturbation τ is imperceptible as soon as $\|\tau\|_p \leq \alpha_p$. It is worth noting that ℓ_p norms are only surrogates for the perception distance, for which it is still an open question to give a formal definition. In this paper, we only focus on robustness on ℓ_p norms. The literature on adversarial attacks for image classification usually uses either an ℓ_∞ norm akin (Madry et al., 2018) or an ℓ_2 norm akin (Carlini & Wagner, 2017) as a surrogate for imperceptibility. Other authors such as Chen et al. (2018) and Papernot et al. (2016) also used an ℓ_1 norm or an ℓ_0 semi-norm.

To account for adversaries possibly manipulating the input images, one needs to revisit the standard risk minimization by incorporating the adversary in the problem. The goal becomes to minimize the *worst-case* risk under α_p -bounded manipulations. We call this problem the *adversarial risk minimization*. It writes

$$\inf_{\mathbf{h} \in \mathcal{H}} \mathcal{R}^{\text{adv}}(\mathbf{h}; \alpha_p) := \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[\sup_{\tau \in B_p(\alpha_p)} \mathcal{L}_{0/1}(\mathbf{h}(x + \tau), y) \right], \quad (3)$$

where $B_p(\alpha_p) := \{\tau \in \mathcal{X} \text{ s.t. } \|\tau\|_p \leq \alpha_p\}$. In this new formulation, the adversary focuses on optimizing the inner maximization, while the learner tries to get the best hypothesis from \mathcal{H} “under attack”. By analogy with the standard setting, given n training examples $\mathcal{S} := \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$, we want to find an approximate solution to the adversarial risk minimization by studying its empirical counterpart, the *empirical adversarial risk minimization*. This optimization problem writes

$$\inf_{\mathbf{h} \in \mathcal{H}} \mathcal{R}_S^{\text{adv}}(\mathbf{h}; \alpha_p) := \frac{1}{n} \sum_{i=1}^n \sup_{\tau \in B_p(\alpha_p)} \mathcal{L}_{0/1}(\mathbf{h}(\mathbf{x}_i + \tau), y_i). \quad (4)$$

In the presence of an adversary, two major issues appear in the empirical risk minimization. First, as recently pointed out by Madry et al. (2018), the adversarial generalization error (i.e. the gap between the empirical adversarial risk and the adversarial risk) can be much larger than in the standard setting. Indeed, the adversary makes the problem

dependent on the dimension of \mathcal{X} . Hence, in high-dimension (e.g. for images) one needs much more samples to classify correctly as pointed out by Schmidt et al. (2018) as well as Simon-Gabriel et al. (2019). Moreover, finding an approximate solution to the adversarial risk minimization is not always sufficient. Indeed, recent works by Tsipras et al. (2019) and Zhang et al. (2019) gave theoretical evidence that training a robust model may lead to an increase of its standard risk. Hence finding a good approximation for Problem (3) may lead to a poor solution for Problem (1). Accordingly, it is natural to wonder whether we can find a class of models \mathcal{H} for which we can control both the standard and adversarial risks?

In this paper, we provide answers to the above question by conducting an in depth analysis of a special class of models called randomized classifiers, i.e. classifiers that output random variables instead of labels. Our main contributions summarize as follows.

1.3 Contributions

Our first contribution consists in studying randomized classifiers. By analogy with the deterministic case, we define a notion of robustness for randomized classifiers. This definition amounts to making the classifier locally Lipschitz with respect to the ℓ_p norm on \mathcal{X} , and a probability metric on \mathcal{Y} (e.g. the total variation distance or the Renyi divergence). More precisely, if we denote D the probability metric at hand, a randomized classifier m is called (α_p, ϵ) -robust w.r.t. D if for any $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$

$$\|\mathbf{x} - \mathbf{x}'\|_p \leq \alpha_p \implies D(m(\mathbf{x}), m(\mathbf{x}')) \leq \epsilon.$$

Denoting $\mathcal{M}_D(\alpha_p, \epsilon)$ the class of randomized classifiers that respect this local Lipschitz condition, we present the following results.

1. If D is either the total variation distance or the Renyi divergence, we show that for any $m \in \mathcal{M}_D(\alpha_p, \epsilon)$, we can upper-bound the gap between the risk and the adversarial risk of m . Notably, if D is the total variation distance, for any $m \in \mathcal{M}_D(\alpha_p, \epsilon)$ we have $\mathcal{R}^{\text{adv}}(m; \alpha_p) - \mathcal{R}(m) \leq \epsilon$. Hence, ϵ controls the maximal trade-off between robust and standard accuracy for locally Lipschitz randomized classifier. We demonstrate similar results when D is the Renyi divergence showing that $\mathcal{R}^{\text{adv}}(m; \alpha_p) - \mathcal{R}(m) \leq 1 - O(e^{-\epsilon})$. This means that, for the class of locally Lipschitz randomized classifiers, solving the risk minimization problem, i.e. Problem (1), gives an approximate solution to the adversarial risk minimization problem, i.e. Problem (3), up to an additive factor that depends on the robustness parameter ϵ .
2. We devise an upper-bound on the generalization gap of any m in $\mathcal{M}_D(\alpha_p, \epsilon)$. In particular, when D is the total variation distance, we demonstrate that for any $m \in \mathcal{M}_D(\alpha_p, \epsilon)$ we have

$$\mathcal{R}(m) - \mathcal{R}_S(m) \leq O\left(\sqrt{\frac{N \times K}{n}}\right) + \epsilon,$$

where N is the external α_p -covering number of the input samples. This means that, when $N/n \xrightarrow{n \rightarrow \infty} 0$, solving the empirical risk minimization problem, i.e. Problem (2), on $\mathcal{M}_D(\alpha_p, \epsilon)$ provides an approximate solution to the risk minimization problem, i.e. Problem (1). Since we can also bound the gap between the adversarial and the standard risk, we can combine the two results to bound the adversarial generalization

gap on $\mathcal{M}_D(\alpha_p, \epsilon)$. Note however, that this result relies on a strong assumption on \mathcal{X} that does not always avoid dimensionality issues. The problem of finding a subclass of $\mathcal{M}_D(\alpha_p, \epsilon)$ that provides tighter generalization bounds is an open question.

For our second contribution, we present a practical way to design this class $\mathcal{M}(\alpha_p, \epsilon)$ by using a simple yet efficient noise injection scheme. This allows us to build randomized classifiers from state-of-the-art machine learning models, including deep neural networks. More precisely our contribution is as follows.

1. Based on information-theoretic properties of the total variation distance and the Renyi divergence (e.g. the data processing inequality) we design a noise injection scheme to turn a state-of-the-art machine learning model into a robust randomized classifier. More formally, let us denote Φ the c.d.f. of a standard Gaussian distribution. Let us consider \mathbf{h} a deterministic hypothesis, we show that the randomized classifier $m : \mathbf{x} \mapsto \mathbf{h}(\mathbf{x} + n)$ with $n \sim \mathcal{N}(0, \sigma^2 I_d)$ is both $(\alpha_2, \frac{(\alpha_2)^2}{2\sigma})$ -robust w.r.t. the Renyi divergence and $(\alpha_2, 2\Phi(\frac{\alpha_2}{2\sigma}) - 1)$ -robust w.r.t. the total variation distance. Our results on randomized classifiers are applicable to a wide range of machine learning models including deep neural networks.
2. We further corroborate our theoretical results with experiments using deep neural networks on standard image datasets, namely CIFAR-10 and CIFAR-100 (Krizhevsky & Hinton, 2009). These models can simultaneously provide accurate prediction (over 0.82 clean accuracy on CIFAR-10) and reasonable robustness against ℓ_2 adversarial examples (0.45 against ℓ_2 adversaries with magnitude 0.5 on CIFAR-10).

2 Related work

Contrary to other notions such as training corruption, a.k.a. poisoning attacks (Kearns & Li, 1993; Kearns et al., 1994), the theoretical study of adversarial robustness is still in its infancy. So far, empirical observations tend to show that (1) adversarial examples on state-of-the-art models are hard to mitigate and (2) robust training methods give poor generalization performances. Some recent works started to study the problem through the lens of learning theory either to understand the links between robustness and accuracy or to provide bounds on the generalization gap of current learning procedures in the adversarial setting.

2.1 Accuracy versus robustness trade-off

A first line of research (Su et al., 2018; Jetley et al., 2018; Tsipras et al., 2019) suggests that designing robust models might be inconsistent with standard accuracy. These works argue with experiments and toy examples that robust and standard classification are two concurrent problems. Following this line, Zhang et al. (2019) observed that the adversarial risk of any hypothesis \mathbf{h} decomposes as follows,

$$\mathcal{R}^{\text{adv}}(\mathbf{h}; \alpha_p) = \mathcal{R}(\mathbf{h}) + \mathcal{R}_{>0}^{\text{adv}}(\mathbf{h}; \alpha_p), \quad (5)$$

where $\mathcal{R}_{>0}^{\text{adv}}(m; \alpha_p)$ is the amount of risk that the adversary gets with *non-null* perturbations. Looking at Eq. (5), we realize that minimizing the adversarial risk is not enough to control standard accuracy, as one could only optimize over the second term. This indicates that

adversarial risk minimization, i.e. Problem (3), is harder to solve than the standard risk minimization, i.e. Problem (1).

While this indicates that both goals may be difficult to achieve simultaneously, Eq. (5), along with the empirical studies from the literature do not highlight any fundamental trade-off between robustness and accuracy. Moreover, no upper-bound on $\mathcal{R}_{>0}^{\text{adv}}(\mathbf{h}; \alpha_p)$ has been demonstrated yet. Hence the questions whether this trade-off exists and can be controlled remain open. In this paper, we provide a rigorous answer to these questions by identifying classes $\mathcal{M}_D(\alpha_p, \epsilon)$ of randomized classifiers for which we can upper bound the trade-off term $\mathcal{R}_{>0}^{\text{adv}}(\mathbf{m}; \alpha_p)$ for any $\mathbf{m} \in \mathcal{M}_D(\alpha_p, \epsilon)$. Hence, we can control the maximum loss of accuracy that the model can suffer in the adversarial setting. It also challenges the intuitions developed by previous works (Su et al., 2018; Jetley et al., 2018; Tsipras et al., 2019) and argues in favor of using randomized mechanisms as a defense against adversarial attacks.

2.2 Studying adversarial generalization

To further compare the hardness of the two problems, a recent line of research began to explore the notion of adversarial generalization gap. In this line, Schmidt et al. (2018) presented some first intuitions by studying a simplified binary classification framework where \mathcal{D} is a mixture of multi-dimensional Gaussian distributions. In this framework the authors show that without attacks, we only need $O(1)$ training samples to have a small generalization gap. But against an ℓ_∞ adversary, we need $O(\sqrt{d})$ training samples instead. In the discussion of their work, the authors present the problem of obtaining similar results without making any assumption about the distribution as an open problem.

This issue was recently studied using the Rademacher complexity by Khim and Loh (2018), Yin et al. (2019) and Awasthi et al. (2020). These papers relate the adversarial generalization error of linear classifiers and one-hidden layer neural networks with the dimension of the problem. They show that the adversarial generalization depends on the dimension of the problem. At a first glance, the difficulty of adversarial generalization seems to contradict previous conclusions on the link between robustness and generalization presented by Xu and Mannor (2012). But, as we will discuss in the sequel, these results assume that the input space \mathcal{X} can be partitioned in $O(1)$ sub-space in which the classification function has small variations. This assumption may not always hold when dealing with high dimensional input spaces (e.g. images) and very sophisticated classification algorithms (e.g. deep neural networks).

Going further, it should be noted that the generalization gap measures only the difference between empirical and theoretical risks. In practice, the empirical adversarial risk is hard to estimate, since we cannot compute the exact solution to the inner maximization problem. The following question therefore remains open: even if we can set up a learning procedure with a controlled generalization gap, can we give guarantees on the standard and adversarial risks? In this paper, we start answering this question by providing techniques that provably offer both small standard risk and reasonable robustness against adversarial examples (see Sect. 1.3 for more details).

2.3 Defense against adversarial examples based on noise injection

Injecting noise into algorithms to improve train time robustness has been used for ages in detection and signal processing tasks (Zozor & Amblard, 1999; Chapeau-Blondeau & Rousseau, 2004; Mitaim & Kosko, 1998; Grandvalet et al., 1997). It has also been

extensively studied in several machine learning and optimization fields, e.g. robust optimization (Ben-Tal et al., 2009) and data augmentation techniques (Perez & Wang, 2017). Concurrently to our work, noise injection techniques have been adopted by the adversarial defense community under the *randomized smoothing* name. The idea of provable defense through noise injection was first proposed by Lecuyer et al. (2019) and refined by Li et al. (2019), Cohen et al. (2019), Salman et al. (2019) and Yang et al. (2020). The rationale behind randomized smoothing is very simple: smooth h after training by convolution with a Gaussian measure to build a more stable classifier. Our work belongs to the same line of research, but the nature of our results is different. Randomized smoothing is an ensemble method that builds a deterministic classifier by smoothing a pre-trained model with a Gaussian kernel. This scheme requires to compute a Monte-Carlo estimation of the smoothed classifier; hence requiring many rounds of evaluations to output a deterministic label. Our method is based on randomization and only requires one evaluation round for inferring a label, making the prediction randomized and computationally efficient. While randomized smoothing focuses on the construction of certified defenses, we study the generalization properties of randomized mechanisms both in the standard and the adversarial setting. Our analysis presents the fundamental properties of randomized defenses, including (but not limited to) randomized smoothing (c.f. Sect. 7).

This paper is an extended version of a work by Pinot et al. (2019). Since then, we considerably consolidated our theoretical results as follows.

1. Pinot et al. (2019) only studied neural networks defended with noise injection techniques, here we study the much more general class of randomized classifiers which includes, but is not limited to neural networks.
2. We provide a much more detailed analysis of our notion of distributional robustness by presenting an in depth analysis based on the Total variation distance that was missing from (Pinot et al., 2019) (Theorems 1, 5 and 7).
3. Pinot et al. did not analyze the generalization of randomized classifiers. Here, we study the generalization of these classifiers according to the notion of robustness they respect (Theorem 5 and Corollary 1).
4. Last but not least, we added an in-depth discussion on the fundamental properties of randomized classifiers, and how they relate to the notion of randomized smoothing (Sect. 7).

3 Definition of risk and robustness for randomized classifiers

In this work, the goal is to analyze how randomized classifiers can solve the problem of classification in the presence of an adversary. Let us start by defining what we mean by randomized classifiers.

Remark 1 (Note on measurability) Through the paper, we assume every spaces \mathcal{Z} to be associated with a σ -algebra denoted $\mathcal{A}(\mathcal{Z})$. Furthermore, we denote $\mathcal{P}(\mathcal{Z})$ the set of probability distributions defined on the measurable space $(\mathcal{Z}, \mathcal{A}(\mathcal{Z}))$. In the following, for simplicity, we refer to $\mathcal{A}(\mathcal{Z})$ only when necessary.

Definition 1 (*Probabilistic mapping*) Let \mathcal{Z} and \mathcal{Z}' be two arbitrary spaces. A *probabilistic mapping* from \mathcal{Z} to \mathcal{Z}' is a mapping $m : \mathcal{Z} \rightarrow \mathcal{P}(\mathcal{Z}')$, where $\mathcal{P}(\mathcal{Z}')$ is the space of

probability measures on \mathcal{Z}' . When $\mathcal{Z} = \mathcal{X}$ and $\mathcal{Z}' = \mathcal{Y}$, m is called a *randomized classifier*. To get a numerical answer for an input \mathbf{x} , we sample $\hat{y} \sim m(\mathbf{x})$.

Any mapping can be considered as a probabilistic mapping, whether it explicitly considers randomization or not. In fact, any deterministic classifier can be considered as a randomized one, since it can be characterized by a Dirac measure. Accordingly, the definition of a randomized classifier is fully general and equally consider classifiers with or without randomization scheme.

3.1 Risk and adversarial risk for randomized classifiers

To analyze this new hypothesis class, we can adapt the concepts of risk and adversarial risk for a randomized classifier. The loss function we use is the natural extension of the 0/1 loss to the randomized regime. Given a randomized classifier m and a sample $(\mathbf{x}, y) \sim \mathcal{D}$ it writes

$$\mathcal{L}_{0/1}(m(\mathbf{x}), y) := \mathbb{E}_{\hat{y} \sim m(\mathbf{x})} [\mathbb{1}\{\hat{y} \neq y\}]. \tag{6}$$

This loss function evaluates the probability of misclassification of m on a data sample $(\mathbf{x}, y) \sim \mathcal{D}$. Accordingly, the risk of m with respect to \mathcal{D} writes

$$\mathcal{R}(m) := \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [\mathcal{L}_{0/1}(m(\mathbf{x}), y)]. \tag{7}$$

Finally, given m and $(\mathbf{x}, y) \sim \mathcal{D}$, the adversary seeks a perturbation $\boldsymbol{\tau} \in B_p(\alpha_p)$ that maximizes the expected error of the classifier on \mathbf{x} (i.e. $\mathbb{E}_{\hat{y} \sim m(\mathbf{x} + \boldsymbol{\tau})} [\mathbb{1}\{\hat{y} \neq y\}]$). Therefore, the adversarial risk of m under α_p -bounded perturbations writes

$$\mathcal{R}^{\text{adv}}(m; \alpha_p) := \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[\sup_{\boldsymbol{\tau} \in B_p(\alpha_p)} \mathcal{L}_{0/1}(m(\mathbf{x} + \boldsymbol{\tau}), y) \right]. \tag{8}$$

By analogy with the deterministic setting, we denote

$$\mathcal{R}_S(m) := \frac{1}{n} \sum_{i=1}^n \mathcal{L}_{0/1}(m(\mathbf{x}_i), y_i), \text{ and} \tag{9}$$

$$\mathcal{R}_S^{\text{adv}}(m; \alpha_p) := \frac{1}{n} \sum_{i=1}^n \sup_{\boldsymbol{\tau} \in B_p(\alpha_p)} \mathcal{L}_{0/1}(m(\mathbf{x}_i + \boldsymbol{\tau}), y_i), \tag{10}$$

the empirical risks of m for a given training sample $\mathcal{S} := \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$.

3.2 Robustness for randomized classifiers

We could define the notion of robustness for a randomized classifier depending on whether it misclassifies any test sample $(\mathbf{x}, y) \sim \mathcal{D}$. But in practice, neither the adversary nor the model provider have access to the ground-truth distribution \mathcal{D} . Furthermore, in real-world scenarios, one wants to check before its deployment that the model is robust. Therefore, it is required for the classifier to be stable on the regions of the space where it already

classifies correctly. Formally a (deterministic) classifier $c : \mathcal{X} \rightarrow \mathcal{Y}$ is called *robust* if for any $(\mathbf{x}, y) \sim \mathcal{D}$ such that $c(\mathbf{x}) = y$, and for any $\boldsymbol{\tau} \in \mathcal{X}$ one has

$$\|\boldsymbol{\tau}\|_p \leq \alpha_p \implies c(\mathbf{x}) = c(\mathbf{x} + \boldsymbol{\tau}). \quad (11)$$

By analogy with this, we define robustness for a randomized classifier below.

Definition 2 (Robustness for a randomized classifier) A randomized classifier $m : \mathcal{X} \rightarrow \mathcal{P}(\mathcal{Y})$ is called (α_p, ϵ) -robust w.r.t. D if for any $\mathbf{x}, \boldsymbol{\tau} \in \mathcal{X}$, one has

$$\|\boldsymbol{\tau}\|_p \leq \alpha_p \implies D(m(\mathbf{x}), m(\mathbf{x} + \boldsymbol{\tau})) \leq \epsilon.$$

Where D is a metric/divergence between two probability measures. Given such a metric/divergence D , we denote $\mathcal{M}_D(\alpha_p, \epsilon)$ the set of all randomized classifiers that are (α_p, ϵ) -robust w.r.t. D .

Note that we did not add the constraint that m classifies well on $(\mathbf{x}, y) \sim \mathcal{D}$, since it is already encompassed in the probability distribution itself. If the two probabilities $m(\mathbf{x})$ and $m(\mathbf{x} + \boldsymbol{\tau})$ are close, and if $m(\mathbf{x})$ outputs y with high probability, then it will be the same for $m(\mathbf{x} + \boldsymbol{\tau})$. This formulation naturally raises the question of the choice of the metric D . Any choice of metric/divergence will instantiate a notion of adversarial robustness, and it should be carefully selected. In the present work, we focus our study on the total variation distance and the Renyi divergence. The question whether these metrics/divergences are more appropriate than others remains open but these two divergences are sufficiently general to cover a wide range of other definitions (see ‘‘Appendix 2’’ for more details). Furthermore, these notions of distance comply with both a theoretical analysis (Sect. 5) and practical considerations (Sect. 8).

3.3 Divergence and probability metrics

Let us now recall the definition of total variation distance and Renyi divergence. Let \mathcal{Z} be an arbitrary space, and ρ, ρ' be two measures in $\mathcal{P}(\mathcal{Z})$.¹ The *total variation distance* between ρ and ρ' is

$$D_{TV}(\rho, \rho') := \sup_{Z \subset \mathcal{A}(\mathcal{Z})} |\rho(Z) - \rho'(Z)|, \quad (12)$$

where $\mathcal{A}(\mathcal{Z})$ is the σ -algebra associated with the set of measures $\mathcal{P}(\mathcal{Z})$. The total variation distance is one of the most commonly used probability metrics. It admits several very simple interpretations, and is a very useful tool in many mathematical fields such as probability theory, Bayesian statistics or optimal transport (Villani, 2003; Robert, 2007; Peyré & Cuturi, 2019). In optimal transport, it can be rewritten as the solution of the Monge-Kantorovich problem with the cost function $\text{cost}(z, z') = \mathbb{1}\{z \neq z'\}$,

$$D_{TV}(\rho, \rho') = \inf \int_{\mathcal{Z}^2} \mathbb{1}\{z \neq z'\} d\pi(z, z'), \quad (13)$$

¹ Recall from Definition 1 that $\mathcal{P}(\mathcal{Z})$ is the set of probability measures on \mathcal{Z} .

where the infimum is taken over all joint probability measures π in $\mathcal{P}(\mathcal{Z} \times \mathcal{Z})$ with marginals ρ and ρ' . According to this interpretation, it seems quite natural to consider the total variation distance as a relaxation of the trivial distance on $[0, 1]$ (for deterministic classifiers).

Let us now suppose that ρ and ρ' admit probability density functions g and g' according to a third measure ν . Then the *Rényi divergence of order β* between ρ and ρ' writes

$$D_\beta(\rho, \rho') := \frac{1}{\beta - 1} \log \int_{\mathcal{Y}} g'(y) \left(\frac{g(y)}{g'(y)} \right)^\beta d\nu(y). \tag{14}$$

The Rényi divergence (Rényi, 1961) is a generalized divergence defined for any β on the interval $[1, \infty]$. It equals the Kullback–Leibler divergence when $\beta \rightarrow 1$, and the maximum divergence when $\beta \rightarrow \infty$. It also has the property of being non-decreasing with respect to β . This divergence is very common in machine learning and Information theory (van Erven & Harremos, 2014), especially in its Kullback-Leibler form as it is widely used as the loss function, i.e. cross entropy, of classification algorithms. In the remaining, we denote $\mathcal{M}_\beta(\alpha_p, \epsilon)$ the set of (α_p, ϵ) -robust classifiers w.r.t. D_β .

Let us now give some properties of these divergences that will be useful for our analysis. First we recall the probability preservation property of the Rényi divergence, first presented by Langlois et al. (2014).

Proposition 1 (Langlois et al., 2014) *Let ρ and ρ' be two measures in $\mathcal{P}(\mathcal{Z})$. Then for any $Z \in \mathcal{A}(\mathcal{Z})$, the following holds,*

$$\rho(Z) \leq (\exp(D_\beta(\rho, \rho')) \rho'(Z))^{\frac{\beta-1}{\beta}}.$$

Now thanks to previous works by Gilardoni (2010) and Vajda (1970), we also get the following results relating the total variation distance and the Rényi divergence.

Proposition 2 (Inequality between total variation and Rényi divergence) *Let ρ and ρ' be two measures in $\mathcal{P}(\mathcal{Z})$, and $\beta \geq 1$. Then the following holds,*

$$D_{TV}(\rho, \rho') \leq \min \left\{ \frac{3}{2} \left(\sqrt{1 + \frac{4D_\beta(\rho, \rho')}{9}} - 1 \right)^{1/2}, \frac{\exp(D_\beta(\rho, \rho') + 1) - 1}{\exp(D_\beta(\rho, \rho') + 1) + 1} \right\}.$$

Proof Thanks to Gilardoni (2010), one has

$$D_1(\rho, \rho') \geq 2D_{TV}(\rho, \rho')^2 + \frac{4D_{TV}(\rho, \rho')^4}{9}.$$

From which it follows that

$$D_{TV}(\rho, \rho') \leq \frac{3}{2} \left(\sqrt{1 + \frac{4D_1(\rho, \rho')}{9}} - 1 \right)^{1/2}.$$

Moreover, using inequality from Vajda (1970), one gets

$$D_1(\rho, \rho') + 1 \geq \log \left(\frac{1 + D_{TV}(\rho, \rho')}{1 - D_{TV}(\rho, \rho')} \right).$$

This inequality leads to the following

$$\frac{\exp(D_1(\rho, \rho') + 1) - 1}{\exp(D_1(\rho, \rho') + 1) + 1} \geq D_{TV}(\rho, \rho').$$

By combining the above inequalities and by monotony of Renyi divergence regarding β , one obtains the expected result. \square

From now on, we denote $\mathcal{M}_{TV}(\alpha, \epsilon)$ and $\mathcal{M}_\beta(\alpha, \epsilon)$ the set of (α, ϵ) -robust classifiers respectively for D_{TV} and D_β . The next section gives bounds on the generalization gap in the standard and the adversarial settings for these specific hypothesis classes.

4 Risks' gap and generalization gap for robust randomized classifiers

As discussed in Sect. 2.1, we can always decompose the adversarial risk of a classifier $\mathcal{R}^{\text{adv}}(\mathbf{m}; \alpha_p)$ in two terms. First the standard risk $\mathcal{R}(\mathbf{m})$ and second the amount of risk the adversary creates with non-zero perturbations $\mathcal{R}_{>0}^{\text{adv}}(\mathbf{m}; \alpha_p)$. Hence minimizing $\mathcal{R}(\mathbf{m})$ can give poor values for $\mathcal{R}^{\text{adv}}(\mathbf{m}; \alpha_p)$ and vice-versa. In this section, we upper-bound the risks' gap $\mathcal{R}_{>0}^{\text{adv}}(\mathbf{m}; \alpha_p)$, i.e. the gap between the risk and the adversarial risk of a robust classifier.

4.1 Risks' gap for robust classifiers w.r.t. D_{TV}

First, let us consider $\mathbf{m} \in \mathcal{M}_{TV}(\alpha_p, \epsilon)$. We can control the loss of accuracy under attack of this classifier with the robustness parameter ϵ .

Theorem 3 (Risk's gap for robust classifiers w.r.t D_{TV}) *Let $\mathbf{m} \in \mathcal{M}_{TV}(\alpha_p, \epsilon)$. Then we have*

$$\mathcal{R}^{\text{adv}}(\mathbf{m}; \alpha_p) \leq \mathcal{R}(\mathbf{m}) + \epsilon.$$

Proof Let \mathbf{m} be an (α_p, ϵ) -robust classifier w.r.t. D_{TV} , $(\mathbf{x}, y) \sim \mathcal{D}$ and $\boldsymbol{\tau} \in \mathcal{X}$ such that $\|\boldsymbol{\tau}\|_p \leq \alpha_p$. By definition of the 0/1 loss we have

$$\mathcal{L}_{0/1}(\mathbf{m}(\mathbf{x} + \boldsymbol{\tau}), y) = \mathbb{E}_{\hat{y} \sim \mathbf{m}(\mathbf{x} + \boldsymbol{\tau})} [\mathbb{1}\{\hat{y} \neq y\}].$$

Furthermore, by definition of the total variation distance we have

$$\mathbb{E}_{\hat{y} \sim \mathbf{m}(\mathbf{x} + \boldsymbol{\tau})} [\mathbb{1}\{\hat{y} \neq y\}] - \mathbb{E}_{\hat{y} \sim \mathbf{m}(\mathbf{x})} [\mathbb{1}\{\hat{y} \neq y\}] \leq D_{TV}(\mathbf{m}(\mathbf{x}), \mathbf{m}(\mathbf{x} + \boldsymbol{\tau})).$$

Since $\mathbf{m} \in \mathcal{M}_{TV}(\alpha_p, \epsilon)$, the above amounts to write

$$\mathcal{L}_{0/1}(\mathbf{m}(\mathbf{x} + \boldsymbol{\tau}), y) - \mathcal{L}_{0/1}(\mathbf{m}(\mathbf{x}), y) \leq \epsilon.$$

Finally, this holds for any $(\mathbf{x}, y) \sim \mathcal{D}$ and any α_p bounded perturbation $\boldsymbol{\tau}$, then we get

$$\mathbb{E}_{(\mathbf{x},y)\sim\mathcal{D}} \left[\sup_{\boldsymbol{\tau}\in B_p(\alpha_p)} \mathcal{L}_{0/1}(\mathbf{m}(\mathbf{x} + \boldsymbol{\tau}), y) \right] - \mathbb{E}_{(\mathbf{x},y)\sim\mathcal{D}} [\mathcal{L}_{0/1}(\mathbf{m}(\mathbf{x}), y)] \leq \epsilon.$$

The above inequality concludes the proof. □

This result means that if we can design a class $\mathcal{M}_{TV}(\alpha_p, \epsilon)$ with small enough ϵ , then minimizing the risk of $m \in \mathcal{M}_{TV}(\alpha_p, \epsilon)$ is also sufficient to control the adversarial risk. It is relatively easy to obtain, but it has an interesting consequence on the understanding we have of the trade-off between robustness and accuracy. It says that there exists some classes of randomized classifiers for which robustness and standard accuracy may not be at odds, since we can upper-bound the maximal loss of accuracy the model may suffer under attack. This questions previous intuitions developed on deterministic classifiers by Su et al. (2018), Jetley et al. (2018), Tsipras et al. (2019) and Zhang et al. (2019) and advocates for the use of randomization schemes as defenses against adversarial attacks. Note, however, that we did not evade the trade-off between robustness and accuracy, we only showed that with certain hypothesis classes it can be controlled.

4.2 Risks’ gap for robust classifiers w.r.t. D_β

We now extend the previous results the Renyi divergence. We show that, for any randomized classifier in $\mathcal{M}_\beta(\alpha_p, \epsilon)$, we can bound the gap between the risk and the adversarial risk of m . Using the Renyi divergence, the factor that controls the classifier’s loss of accuracy under attack can be either multiplicative or additive, and depends both on the robustness parameter ϵ and on the divergence parameter β .

Theorem 4 (Multiplicative risks’ gap for Renyi-robust classifiers) *Let $m \in \mathcal{M}_\beta(\alpha_p, \epsilon)$. Then we have*

$$\mathcal{R}^{\text{adv}}(m; \alpha_p) \leq (e^\epsilon \mathcal{R}(m))^{\frac{\beta-1}{\beta}}.$$

Proof Let m be an (α_p, ϵ) -robust classifier w.r.t. D_β , $(\mathbf{x}, y) \sim \mathcal{D}$ and $\boldsymbol{\tau} \in \mathcal{X}$ such that $\|\boldsymbol{\tau}\|_p \leq \alpha_p$. With the same reasoning as above, and with Proposition 1, we get

$$\begin{aligned} \mathcal{L}_{0/1}(\mathbf{m}(\mathbf{x} + \boldsymbol{\tau}), y) &= \mathbb{E}_{\hat{y}\sim m(\mathbf{x}+\boldsymbol{\tau})} [\mathbb{1}\{\hat{y} \neq y\}] \\ &= \mathbb{P}_{\hat{y}\sim m(\mathbf{x}+\boldsymbol{\tau})} [\hat{y} \neq y] \\ &\leq (e^{D_\beta(m(\mathbf{x}+\boldsymbol{\tau}), m(\mathbf{x}))} \mathbb{P}_{\hat{y}\sim m(\mathbf{x})} [\hat{y} \neq y])^{\frac{\beta-1}{\beta}} \quad (\text{Prop. 1}) \\ &= (e^{D_\beta(m(\mathbf{x}+\boldsymbol{\tau}), m(\mathbf{x}))} \mathbb{E}_{\hat{y}\sim m(\mathbf{x})} [\mathbb{1}\{\hat{y} \neq y\}])^{\frac{\beta-1}{\beta}} \\ &\leq (e^\epsilon \mathcal{L}_{0/1}(\mathbf{m}(\mathbf{x}), y))^{\frac{\beta-1}{\beta}}. \end{aligned}$$

Since this holds for any $(\mathbf{x}, y) \sim \mathcal{D}$ and any α_p bounded perturbation $\boldsymbol{\tau}$, we get

$$\begin{aligned} \mathcal{R}^{\text{adv}}(\mathbf{m}; \alpha_p) &= \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[\sup_{\boldsymbol{\tau} \in B_p(\alpha_p)} \mathcal{L}_{0/1}(\mathbf{m}(\mathbf{x} + \boldsymbol{\tau}), y) \right] \\ &\leq \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[e^{\frac{\beta-1}{\beta}\epsilon} \mathcal{L}_{0/1}(\mathbf{m}(\mathbf{x}), y)^{\frac{\beta-1}{\beta}} \right] \\ &\leq e^{\frac{\beta-1}{\beta}\epsilon} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[\mathcal{L}_{0/1}(\mathbf{m}(\mathbf{x}), y)^{\frac{\beta-1}{\beta}} \right]. \end{aligned}$$

Finally, using the Jensen inequality, one gets

$$\leq e^{\frac{\beta-1}{\beta}\epsilon} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[\mathcal{L}_{0/1}(\mathbf{m}(\mathbf{x}), y) \right]^{\frac{\beta-1}{\beta}} = (e^\epsilon \mathcal{R}(\mathbf{m}))^{\frac{\beta-1}{\beta}}.$$

The above inequality concludes the proof. □

This first result gives a multiplicative bound on the gap between the standard and adversarial risks. This means that if we can design a class $\mathcal{M}_\beta(\alpha_p, \epsilon)$ with small enough ϵ , and big enough β , then minimizing the risk of any $\mathbf{m} \in \mathcal{M}_\beta(\alpha_p, \epsilon)$ is sufficient to also minimize the adversarial risk of \mathbf{m} . Nevertheless, multiplicative factors are not easy to analyze.

Remark 2 More general bounds can be computed if we assume that for every randomized classifier \mathbf{m} there exists a convex function \mathbf{f} such that for all \mathbf{x} and $\boldsymbol{\tau}$ with $\|\boldsymbol{\tau}\|_p \leq \alpha_p$, we have $\mathbf{m}(\mathbf{x})(Z) \leq \mathbf{f}(\mathbf{m}(\mathbf{x} + \boldsymbol{\tau})(Z))$ for all measurable sets Z . In this case, we get $\mathcal{R}^{\text{adv}}(\mathbf{m}; \alpha_p) \leq \mathbf{f}(\mathcal{R}(\mathbf{m}))$. This has a close link with randomized smoothing (Cohen et al., 2019) and f -differential privacy (Bu et al., 2020) where both try to fit the best possible \mathbf{f} using Neyman–Pearson lemma.

The following result provides an additive counterpart to Theorem 4. It gives a control over the loss of accuracy under attack with respect to the robustness parameter ϵ and the Shannon entropy of \mathbf{m} .

Theorem 5 (Additive risks’ gap for Renyi-robust classifiers) *Let $\mathbf{m} \in \mathcal{M}_\beta(\alpha_p, \epsilon)$, then we have*

$$\mathcal{R}^{\text{adv}}(\mathbf{m}; \alpha_p) - \mathcal{R}(\mathbf{m}) \leq 1 - e^{-\epsilon} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{|\mathcal{X}}} \left[e^{-H(\mathbf{m}(\mathbf{x}))} \right]$$

where H is the Shannon entropy (i.e. for any $\rho \in \mathcal{P}(\mathcal{Y})$, $H(\rho) = -\sum_{k \in \mathcal{Y}} \rho_k \log(\rho_k)$) and $\mathcal{D}_{|\mathcal{X}}$ is the marginal distribution of \mathcal{D} for \mathcal{X} .

Proof Let $\mathbf{m} \in \mathcal{M}_\beta(\alpha_p, \epsilon)$, then

$$\begin{aligned} &\mathcal{R}^{\text{adv}}(\mathbf{m}; \alpha_p) - \mathcal{R}(\mathbf{m}) \\ &= \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[\sup_{\boldsymbol{\tau} \in B_p(\alpha_p)} \mathcal{L}_{0/1}(\mathbf{m}(\mathbf{x} + \boldsymbol{\tau}), y) - \mathcal{L}_{0/1}(\mathbf{m}(\mathbf{x}), y) \right]. \end{aligned}$$

By definition of the 0/1 loss, this amounts to write

$$\begin{aligned}
 &= \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[\sup_{\tau \in B_p(\alpha_p)} \mathbb{E}_{\hat{y}_{\text{adv}} \sim m(x+\tau), \hat{y} \sim m(x)} \left[\mathbb{1}(\hat{y}_{\text{adv}} \neq y) - \mathbb{1}(\hat{y} \neq y) \right] \right] \\
 &\leq \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[\sup_{\tau \in B_p(\alpha_p)} \mathbb{E}_{\hat{y}_{\text{adv}} \sim m(x+\tau), \hat{y} \sim m(x)} \left[\mathbb{1}(\hat{y}_{\text{adv}} \neq \hat{y}) \right] \right] \\
 &= \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[\sup_{\tau \in B_p(\alpha_p)} \mathbb{P}_{\hat{y}_{\text{adv}} \sim m(x+\tau), \hat{y} \sim m(x)} \left[\hat{y}_{\text{adv}} \neq \hat{y} \right] \right] \\
 &= \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[\sup_{\tau \in B_p(\alpha_p)} 1 - \mathbb{P}_{\hat{y}_{\text{adv}} \sim m(x+\tau), \hat{y} \sim m(x)} \left[\hat{y}_{\text{adv}} = \hat{y} \right] \right] \\
 &= \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[\sup_{\tau \in B_p(\alpha_p)} 1 - \sum_{i=1}^K m(x)_i \times m(x + \tau)_i \right].
 \end{aligned}$$

Now, note that for any $(x, y) \sim \mathcal{D}$ and $\tau \in \mathcal{X}$, by definition of a probability vector in $\mathcal{P}(\mathcal{Y})$, and thanks to Jensen inequality we can write

$$\sum_{i=1}^K m(x)_i \times m(x + \tau)_i \geq \exp \left(\sum_{i=1}^K m(x)_i \log m(x + \tau)_i \right).$$

Then by definition of the entropy and the Kullback Leibler divergence we have

$$\exp \left(\sum_{i=1}^K m(x)_i \log m(x + \tau)_i \right) = \exp \left(-D_1(m(x), m(x + \tau)) - H(m(x)) \right).$$

Finally, by combining the above inequalities and since $m \in \mathcal{M}_\beta(\alpha_p, \epsilon)$ we get

$$\begin{aligned}
 &\mathbb{E}_{(x,y) \sim \mathcal{D}} \left[\sup_{\tau \in B_p(\alpha_p)} \mathbb{P}_{\hat{y}_{\text{adv}} \sim m(x+\tau), \hat{y} \sim m(x)} \left(\hat{y}_{\text{adv}} \neq \hat{y} \right) \right] \\
 &\leq \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[\sup_{\tau \in B_p(\alpha_p)} 1 - e^{-D_1(m(x), m(x+\tau)) - H(m(x))} \right] \\
 &\leq \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[1 - e^{-\epsilon - H(m(x))} \right] = 1 - e^{-\epsilon} \mathbb{E}_{x \sim \mathcal{D}_X} \left[e^{-H(m(x))} \right].
 \end{aligned}$$

The above inequality concludes the proof. □

This result is interesting because it relates the accuracy of m with the bound we obtain. In words, when $m(x)$ has large entropy (i.e. $H(m(x)) \rightarrow \log(K)$) the output distribution tends towards the uniform distribution; hence $\epsilon \rightarrow 0$. This means that the classifier is very robust but also completely inaccurate, since it outputs classes uniformly at random. On the opposite, if $H(m(x)) \rightarrow 0$, then $\epsilon \rightarrow \infty$. The classifier may be accurate, but it is not robust anymore (at least according to our definition). Hence we need to find a classifier that achieves a trade-off between robustness and accuracy.

5 Standard generalization gap

In this section we devise generalization gap bounds for randomized classifiers when they are robust according either to the total variation distance or the Renyi divergence. To do so, we upper-bound the Rademacher complexity of the loss space for TV-robust classifiers

$$\mathcal{L}_{\mathcal{M}_{TV}(\alpha_p, \epsilon)} := \{(\mathbf{x}, y) \mapsto \mathcal{L}_{0/1}(\mathbf{h}(\mathbf{x}), y) \mid \mathbf{m} \in \mathcal{M}_{TV}(\alpha_p, \epsilon)\}.$$

The *empirical Rademacher complexity*, first introduced by Bartlett and Mendelson (2002), is one of the standard measures of generalization gap. It is particularly useful to obtain quality bounds for complex classes such as neural networks since it does not depend on the number of parameters in the network contrary to combinatorial notions such as the *VC dimension*.

Definition 3 (*Rademacher complexity*) For any class of real-valued functions $\mathcal{F} := \{(\mathbf{x}, y) \mapsto \mathbb{R}\}$, given a training sample $\mathcal{S} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$, the *empirical Rademacher complexity* of \mathcal{F} is defined as

$$\mathfrak{R}_{\mathcal{S}}(\mathcal{F}) := \frac{1}{n} \mathbb{E}_{r_i} \left[\sup_{f \in \mathcal{F}} \sum_{i=1}^n r_i f(\mathbf{x}_i, y_i) \right],$$

with r_i i.i.d. drawn from a Rademacher measure, i.e. $\mathbb{P}(r_i = 1) = \mathbb{P}(r_i = -1) = \frac{1}{2}$.

The empirical Rademacher complexity measures the uniform convergence rate of the empirical risk towards the risk on the function class \mathcal{F} as demonstrated by Mohri et al. (2018). Thanks to this notion of complexity, we can bound with high probability the generalization gap of any hypothesis \mathbf{m} in a class \mathcal{M} .

Theorem 6 (Mohri et al., 2018) *Let \mathcal{M} be a class of possibly randomized classifiers and $\mathcal{L}_{\mathcal{M}} := \{\mathcal{L}_{\mathbf{m}} : (\mathbf{x}, y) \mapsto \mathcal{L}_{0/1}(\mathbf{m}(\mathbf{x}), y) \mid \mathbf{m} \in \mathcal{M}\}$. Then for any $\delta \in (0, 1)$, with probability at least $1 - \delta$, the following holds for any $\mathbf{m} \in \mathcal{M}_{TV}(\alpha_p, \epsilon)$,*

$$\mathcal{R}(\mathbf{m}) - \mathcal{R}_{\mathcal{S}}(\mathbf{m}) \leq 2\mathfrak{R}_{\mathcal{S}}(\mathcal{L}_{\mathcal{M}}) + 3\sqrt{\frac{\ln(2/\delta)}{2n}}.$$

5.1 Generalization error for robust classifiers

Accordingly, we want to upper bound the empirical Rademacher complexity of $\mathcal{L}_{\mathcal{M}_{TV}(\alpha_p, \epsilon)}$, which motivates the following definition.

Definition 4 (α -covering and external covering number) Let us consider $(\mathcal{X}, \|\cdot\|_p)$ a vector space equipped with the ℓ_p norm, $B \subset \mathcal{X}$ and $\alpha \geq 0$. Then

- $C = \{\mathbf{c}_1, \dots, \mathbf{c}_m\}$ is an α -covering of B for the ℓ_p norm if for any $\mathbf{x} \in B$ there exists $\mathbf{c}_i \in C$ such that $\|\mathbf{x} - \mathbf{c}_i\|_p \leq \alpha$.

- The external covering number of B writes $N(B, \|\cdot\|_p, \alpha)$. It is the minimal number of points one needs to build an α -covering of B for the ℓ_p norm.

The covering number is a well-known measure that is often used in statistical learning theory (Shalev-Shwartz & Ben-David, 2014) and asymptotic statistics (Van der Vaart, 2000) to evaluate the complexity of a set of functions. Here we use it to evaluate the number of ℓ_p balls we need to cover the training samples, which gives us the following bound on the Rademacher complexity of $\mathcal{L}_{\mathcal{M}_{TV}(\alpha_p, \epsilon)}$.

Theorem 7 (Rademacher complexity for TV-robust classifiers) *Let $\mathcal{L}_{\mathcal{M}_{TV}(\alpha_p, \epsilon)}$ be the loss function class associated with $\mathcal{M}_{TV}(\alpha_p, \epsilon)$. Then, for any $\mathcal{S} := \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$, the following holds,*

$$\mathfrak{R}_{\mathcal{S}}\left(\mathcal{L}_{\mathcal{M}_{TV}(\alpha_p, \epsilon)}\right) \leq \sqrt{\frac{N \times K}{n}} + \epsilon.$$

where $N = N(\{\mathbf{x}_1, \dots, \mathbf{x}_n\}, \|\cdot\|_p, \alpha_p)$ is the α_p -external covering number of the inputs $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ for the ℓ_p norm and $K = |\mathcal{Y}|$ is the number of labels in the classification task.

Proof We denote $\mathcal{S} := \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ and $N = N(\{\mathbf{x}_1, \dots, \mathbf{x}_n\}, \|\cdot\|_p, \alpha_p)$. By definition of a covering number, there exists $C = \{c_1, \dots, c_N\}$ an α_p -covering of $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ for the ℓ_p norm. Furthermore, for $j \in \{1, \dots, N\}$ and $y \in \{1, \dots, K\}$, we define

$$E_{y,j} = \left\{ i \in \{1, \dots, n\} \text{ s.t. } y_i = y \text{ and } \underset{l \in \{1, \dots, N\}}{\operatorname{argmin}} \|x_i - c_l\| = j \right\}.$$

We also denote $E_j = \cup_{y \in [K]} E_{y,j}$. Finally, we denote $\mathcal{L}_m : (\mathbf{x}, y) \mapsto \mathcal{L}_{0/1}(m(\mathbf{x}), y)$. Then, by definition of the empirical Rademacher complexity, we can write

$$\mathfrak{R}_{\mathcal{S}}\left(\mathcal{L}_{\mathcal{M}_{TV}(\alpha_p, \epsilon)}\right) = \frac{1}{n} \mathbb{E}_{r_i} \left[\sup_{m \in \mathcal{M}_{TV}(\alpha_p, \epsilon)} \sum_{i=1}^n r_i \mathcal{L}_m(\mathbf{x}_i, y_i) \right].$$

Then we can use E_j to write

$$\mathfrak{R}_{\mathcal{S}}\left(\mathcal{L}_{\mathcal{M}_{TV}(\alpha_p, \epsilon)}\right) = \frac{1}{n} \mathbb{E}_{r_i} \left[\sup_{m \in \mathcal{M}_{TV}(\alpha_p, \epsilon)} \sum_{j=1}^N \sum_{i \in E_j} r_i \mathcal{L}_m(\mathbf{x}_i, y_i) \right].$$

Furthermore for any $m \in \mathcal{M}_{TV}(\alpha_p, \epsilon)$ and $i \in E_j$, there exists $\epsilon_i \in [-\epsilon, \epsilon]$ such that: $\mathcal{L}_m(\mathbf{x}_i, y_i) = \mathcal{L}_m(c_j, y_i) + \epsilon_i$. Then we have

$$\begin{aligned} \mathfrak{R}_{\mathcal{S}}\left(\mathcal{L}_{\mathcal{M}_{TV}(\alpha_p, \epsilon)}\right) &\leq \frac{1}{n} \mathbb{E}_{r_i} \left[\sup_{m \in \mathcal{M}_{TV}(\alpha_p, \epsilon)} \sum_{j=1}^N \sum_{i \in E_j} r_i \mathcal{L}_m(c_j, y_i) \right] \\ &\quad + \frac{1}{n} \mathbb{E}_{r_i} \left[\sup_{\epsilon_i \in [-\epsilon, \epsilon]} \sum_{j=1}^N \sum_{i \in E_j} r_i \epsilon_i \right]. \end{aligned}$$

Let us start by studying the second term. We have

$$\frac{1}{n} \mathbb{E}_{r_i} \left[\sup_{\epsilon_i \in [-\epsilon, \epsilon]} \sum_{j=1}^N \sum_{i \in E_j} r_i \epsilon_i \right] = \frac{1}{n} \mathbb{E}_{r_i} \left[\sup_{\epsilon_i \in [-\epsilon, \epsilon]} \sum_{i=1}^n r_i \epsilon_i \right] = \frac{1}{n} \sum_{i=1}^n \epsilon = \epsilon.$$

Now looking at the first term. Since $\mathcal{L}_m(\mathbf{x}, y) \in [0, 1]$ for all (\mathbf{x}, y) we have

$$\begin{aligned} & \frac{1}{n} \mathbb{E}_{r_i} \left[\sup_{m \in \mathcal{M}_{TV}(\alpha_p, \epsilon)} \sum_{j=1}^N \sum_{i \in E_j} r_i \mathcal{L}_m(\mathbf{c}_j, y_i) \right] \\ &= \frac{1}{n} \mathbb{E}_{r_i} \left[\sup_{m \in \mathcal{M}_{TV}(\alpha_p, \epsilon)} \sum_{j=1}^N \sum_{y=1}^K \mathcal{L}_m(\mathbf{c}_j, y) \sum_{i \in E_{y,j}} r_i \right] \\ &\leq \frac{1}{n} \mathbb{E}_{r_i} \left[\sum_{j=1}^N \sum_{y=1}^K \left| \sum_{i \in E_{y,j}} r_i \right| \right]. \end{aligned}$$

Finally using the Khintchine inequality and the Cauchy Schartz inequality we get

$$\begin{aligned} \frac{1}{n} \mathbb{E}_{r_i} \left[\sum_{j=1}^N \sum_{y=1}^K \left| \sum_{i \in E_{y,j}} r_i \right| \right] &\leq \frac{1}{n} \sum_{j=1}^N \sum_{y=1}^K \sqrt{|E_{y,j}|} \quad (\text{Khintchine}) \\ &\leq \frac{1}{n} \sqrt{N \times K} \sqrt{\sum_{j=1}^N \sum_{y=1}^K |E_{y,j}|} \quad (\text{Cauchy}) \\ &= \sqrt{\frac{N \times K}{n}}. \end{aligned}$$

By combining the upper-bounds we have for each term, we get the expected result,

$$\mathfrak{R}_S(\mathcal{L}_{\mathcal{M}_{TV}(\alpha_p, \epsilon)}) \leq \sqrt{\frac{N \times K}{n}} + \epsilon.$$

□

Remark 3 Usually, generalization bounds are involving covering numbers on the hypothesis space using Dudley's entropy integral (Shalev-Shwartz & Ben-David, 2014). In the proposed bound in previous Theorem, it is worth noting that the involved covering number is on the hypothesis space of TV-robust classifiers. This makes a fundamental different between these bounds. Some works (Xu & Mannor, 2012; Petzka et al., 2021) proposed to study the generalization of slowly varying classifiers. The bound they derive are similar to ours, even though they do not apply to the same objects.

The above result means that, if we can cover the n training samples with $O(1)$ balls, then we can bound the generalization gap of any randomized classifier $m \in \mathcal{M}_{TV}(\alpha_p, \epsilon)$ by $O\left(\frac{1}{\sqrt{n}}\right) + \epsilon$. Furthermore, a natural corollary of Theorem 7 bounds the Rademacher complexity of the class $\mathcal{L}_{\mathcal{M}_\beta(\alpha_p, \epsilon)}$.

Corollary 1 Let $\mathcal{L}_{\mathcal{M}_\beta(\alpha_p, \epsilon)}$ be the loss function class associated with $\mathcal{M}_\beta(\alpha_p, \epsilon)$. Then, for any $\mathcal{S} := \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$, the following holds,

$$\mathfrak{R}_{\mathcal{S}}\left(\mathcal{L}_{\mathcal{M}_\beta(\alpha_p, \epsilon)}\right) \leq \sqrt{\frac{N \times K}{n}} + \min\left(\frac{3}{2}\left(\sqrt{1 + \frac{4\epsilon}{9}} - 1\right)^{1/2}, \frac{e^{\epsilon+1} - 1}{e^{\epsilon+1} + 1}\right).$$

where $N = N(\{\mathbf{x}_1, \dots, \mathbf{x}_n\}, \|\cdot\|_p, \alpha_p)$ is the α_p -external covering number of the inputs $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ for the ℓ_p norm.

Proof This corollary is an immediate consequence of Theorem 7 and Proposition 2. □

Thanks to Theorems 6 and 7 and Corollary 1, one can easily bound the generalization gap of robust randomized classifiers.

5.2 Discussion and dimensionality issues

Xu and Mannor (2012) previously studied generalization bounds for learning algorithms based on their robustness. Although we use very different proof techniques, their results and ours are similar. More precisely, both analyses conclude that robust models generalize well if the training samples have a small covering number. Note, however, that we base our formulation on an *adaptive partition* of the samples, while the initial paper from Xu and Mannor (2012) only focuses on a fixed partition of the input space. We refer the reader to the discussion section in Xu and Mannor (2012) for more details.

These findings seem to contradict the current line of works on the hardness of generalization in the adversarial setting. In fact, if the ground truth distribution is sufficiently concentrated (e.g. lies in a low dimensional subspace of \mathbf{x}), a small number of balls can cover \mathcal{S} with high probability; hence $N = O(1)$. This means that we can learn robust classifiers with the same sample complexity as in the standard setting. But if the ground truth distribution is not concentrated enough, the training samples will be far one from another; hence forcing the covering number to be large. In the worse case scenario, we need to cover the whole space $[0, 1]^d$ giving a covering number $N = O\left(\frac{1}{(\alpha_p)^d}\right)$ which is exponential in the dimension of the problem.

Therefore, in the worst-case scenario, our bound is in $O\left(\frac{1}{(\alpha_p)^d \sqrt{n}}\right) + \epsilon$. When α_p is small and the dimension of the problem is high, this bound is too large to give any meaningful insight on the generalization gap of the problem. Therefore, we still need to tighten our analysis to show that robust learning for randomized classifiers is possible in high dimensional spaces.

Remark 4 Note that, we provided a very general result for randomized classifiers under the only assumption that they are robust w.r.t. the total variation distance. Our result applies to any class of classifiers and not only linear classifiers or one-hidden layer neural networks. To build a finer analysis, and to evade the curse of dimensionality, we should consider designing specific sub-classes $\mathcal{M} \subset \mathcal{M}_{TV}(\alpha_p, \epsilon)$ and adapt the proofs to make the term N smaller in the worst-case scenario.

6 Building robust randomized classifiers

In this section we present a simple yet efficient way to transform a non-robust, non-randomized classifier into a robust randomized classifier. To do so, we use a key property of both the Renyi divergence and the total variation distance called the *Data processing inequality*. It is a well-known result from information theory which states that “*post-processing cannot increase information*”. The data processing inequality is as follows.

Theorem 8 (Cover & Thomas, 2012) *Let us consider two arbitrary spaces $\mathcal{Z}, \mathcal{Z}'$, $\rho, \rho' \in \mathcal{P}(\mathcal{Z})$ and $D \in \{D_{TV}, D_\beta\}$. Then for any $\psi : \mathcal{Z} \rightarrow \mathcal{Z}'$ we have*

$$D(\psi\#\rho, \psi\#\rho') \leq D(\rho, \rho'),$$

where $\psi\#\rho$ denotes the pushforward of distribution ρ by ψ .

In the context of robustness to adversarial examples, we use the data processing inequality to ease the design of robust randomized classifiers. In particular, let us suppose that we can build a randomized pre-processing $\mathfrak{p} : \mathcal{X} \rightarrow \mathcal{P}(\mathcal{X})$ such that for any $\mathbf{x} \in \mathcal{X}$ and any α_p -bounded perturbation τ , we have

$$D(\mathfrak{p}(\mathbf{x}), \mathfrak{p}(\mathbf{x} + \tau)) \leq \epsilon, \text{ with } D \in \{D_{TV}, D_\beta\}. \quad (15)$$

Then, thanks to the data processing inequality, we can take any deterministic classifier \mathbf{h} to build an (α_p, ϵ) robust classifier w.r.t D defined as $m : \mathbf{x} \mapsto \mathbf{h}\#\mathfrak{p}(\mathbf{x})$. This considerably simplifies the problem of building a class of robust models. Therefore, we want to build \mathfrak{p} a randomized pre-processing for which we can control the Renyi divergence and/or total variation distance between two inputs. To do this, we analyze the simple procedure of injecting random noise directly on the image before sending it to a classifier. Since the Renyi divergence and the total variation distances are particularly well suited to the study of Gaussian distributions, we first use this type of noise injection. More precisely, in this section, we focus on a mapping that writes as follows.

$$\mathfrak{p} : \mathbf{x} \mapsto \mathcal{N}(\mathbf{x}, \Sigma), \quad (16)$$

for some given non-degenerate covariance matrix $\Sigma \in \mathcal{M}_{d \times d}(\mathbb{R})$. We refer the interested reader to Pinot et al. (2019) for more general classes of noise, namely exponential families. Let us now evaluate the maximal variation of Gaussian pre-processing \mathfrak{p} when applied to an image $\mathbf{x} \in \mathcal{X}$ with and without perturbation.

Lemma 1 *Let $\beta > 1$, $\mathbf{x}, \tau \in \mathcal{X}$ and $\Sigma \in \mathcal{M}_{d \times d}(\mathbb{R})$ a non-degenerate covariance matrix. Let $\rho = \mathcal{N}(\mathbf{x}, \Sigma)$ and $\rho' = \mathcal{N}(\mathbf{x} + \tau, \Sigma)$, then $D_\beta(\rho, \rho') = \frac{\beta}{2} \|\tau\|_{\Sigma^{-1}}^2$.*

Thanks to the above lemma, we know how to evaluate the level of Renyi-robustness that a Gaussian noise pre-processing brings to a classifier. Now that we have this result, thanks to Proposition 2, we can also upper-bound the total variation distance between $\mathcal{N}(\mathbf{x}, \Sigma)$ and $\mathcal{N}(\mathbf{x} + \tau, \Sigma)$. But this bound is not always tight. Besides, we can directly evaluate the total variation distance between two Gaussian distributions as follows.

Lemma 2 Let $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$ and $\Sigma \in \mathcal{M}_{d \times d}(\mathbb{R})$ a non-degenerate covariance matrix. Let $\rho = \mathcal{N}(\mathbf{x}, \Sigma)$ and $\rho' = \mathcal{N}(\mathbf{x} + \boldsymbol{\tau}, \Sigma)$, then $D_{TV}(\rho, \rho') = 2\Phi\left(\frac{\|\boldsymbol{\tau}\|_{\Sigma^{-1}}}{2}\right) - 1$ with Φ the cumulative density function of the standard Gaussian distribution.

Note that both bounds increase with the Mahalanobis norm of $\boldsymbol{\tau}$. Furthermore, we see that the greater the entropy of the Gaussian noise we inject, the smaller the distance between distributions. If we simplify the covariance matrix by setting $\Sigma = \sigma^2 I_d$, it means that we can build more or less robust randomized classifiers against ℓ_2 adversaries, depending on σ .

Theorem 9 (Robustness of Gaussian pre-processing) Let us consider $c : \mathcal{X} \rightarrow \mathcal{Y}$ a deterministic classifier, $\sigma > 0$ and $\mathbf{p} : \mathbf{x} \mapsto \mathcal{N}(\mathbf{x}, \sigma^2 I_d)$ a pre-processing probabilistic mapping. Then the randomized classifier $m := c \# \mathbf{p}$ is

- $(\alpha_2, \frac{(\alpha_2)^2 \beta}{2\sigma})$ -robust w.r.t. D_β against ℓ_2 adversaries.
- $(\alpha_2, 2\Phi\left(\frac{\alpha_2}{2\sigma}\right) - 1)$ -robust w.r.t. D_{TV} against ℓ_2 adversaries.

Proof Let $\mathbf{x}, \boldsymbol{\tau} \in \mathcal{X}$ such that $\|\boldsymbol{\tau}\|_2 \leq \alpha_2$. Thanks to Lemma 1 we have

$$D_\beta(\mathbf{p}(\mathbf{x}), \mathbf{p}(\mathbf{x} + \boldsymbol{\tau})) = \frac{\beta}{2} \|\boldsymbol{\tau}\|_{\Sigma^{-1}}^2 = \frac{\beta}{2\sigma^2} \|\boldsymbol{\tau}\|_2^2 \leq \frac{\beta(\alpha_2)^2}{2\sigma^2}.$$

Similarly, thanks to Lemma 2, we get

$$D_{TV}(\mathbf{p}(\mathbf{x}), \mathbf{p}(\mathbf{x} + \boldsymbol{\tau})) = 2\Phi\left(\frac{\|\boldsymbol{\tau}\|_{\Sigma^{-1}}}{2}\right) - 1 \leq 2\Phi\left(\frac{\alpha_2}{2\sigma}\right) - 1.$$

Finally, from the data processing inequality, i.e. Theorem 8, we get both

$$D_\beta(m(\mathbf{x}), m(\mathbf{x} + \boldsymbol{\tau})) \leq \frac{\beta(\alpha_2)^2}{2\sigma^2},$$

and

$$D_{TV}(m(\mathbf{x}), m(\mathbf{x} + \boldsymbol{\tau})) \leq 2\Phi\left(\frac{\alpha_2}{2\sigma}\right) - 1.$$

The above inequalities conclude the proof. □

Theorem 9 means that we can build simple noise injection schemes as pre-processing of state-of-the-art image classification models and keep track of the maximal loss of accuracy under attack of the resulting randomized classifier. These results also highlight the profound link between randomized classifiers and randomized smoothing as presented by Cohen et al. (2019). Even though our findings are of different nature, both techniques use the same base mechanism (Gaussian noise injection). Therefore, Gaussian pre-processing is a principled defense method that can be analyzed through several standpoints, including certified robustness and statistical learning theory.

7 Discussion: mode preservation property and randomized smoothing

Even though randomized classifiers have some interesting properties regarding generalization error, we can also study them through the prism of deterministic robustness. Let us for example consider the classifier that outputs the class with the highest probability for $\mathbf{m}(\mathbf{x})$, a.k.a. the mode of $\mathbf{m}(\mathbf{x})$. It writes

$$\mathbf{h}_{\text{rob}} : \mathbf{x} \mapsto \operatorname{argmax}_{k \in [K]} \mathbf{m}(\mathbf{x})_k \quad (17)$$

Then checking whether \mathbf{h}_{rob} is robust boils down to demonstrating that the mode of $\mathbf{m}(\mathbf{x})$ does not change under perturbation. It turns out that D_{TV} robust classifiers have this property. We call it the mode preservation property of $\mathcal{M}_{TV}(\alpha_p, \epsilon)$.

Proposition 10 (Mode preservation for D_{TV} -robust classifiers) *Let $\mathbf{m} \in \mathcal{M}_{TV}(\alpha_p, \epsilon)$ be a robust randomized classifier and $\mathbf{x} \in \mathcal{X}$ such that $\mathbf{m}(\mathbf{x})_{(1)} \geq \mathbf{m}(\mathbf{x})_{(2)} + 2\epsilon$. Then, for any $\boldsymbol{\tau} \in \mathcal{X}$, the following holds,*

$$\|\boldsymbol{\tau}\|_p \leq \alpha_p \implies \mathbf{h}_{\text{rob}}(\mathbf{x}) = \mathbf{h}_{\text{rob}}(\mathbf{x} + \boldsymbol{\tau}).$$

Proof Let $\mathbf{x}, \boldsymbol{\tau} \in \mathcal{X}$ such that $\|\boldsymbol{\tau}\|_p \leq \alpha_p$ and $\mathbf{m} \in \mathcal{M}_{TV}(\alpha_p, \epsilon)$ such that

$$\mathbf{m}(\mathbf{x})_{(1)} \geq \mathbf{m}(\mathbf{x})_{(2)} + 2\epsilon.$$

By definition of $\mathcal{M}_{TV}(\alpha_p, \epsilon)$, we have that

$$D_{TV}(\mathbf{m}(\mathbf{x}), \mathbf{m}(\mathbf{x} + \boldsymbol{\tau})) \leq \epsilon.$$

Then, for all $k \in \{1, \dots, K\}$ we have

$$\mathbf{m}(\mathbf{x})_k - \epsilon \leq \mathbf{m}(\mathbf{x} + \boldsymbol{\tau})_k \leq \mathbf{m}(\mathbf{x})_k + \epsilon.$$

Let us denote k^* the index of the biggest value in $\mathbf{m}(\mathbf{x})$, i.e. $\mathbf{m}(\mathbf{x})_{k^*} = \mathbf{m}(\mathbf{x})_{(1)}$. For any $k \in \{1, \dots, K\}$ with $k \neq k^*$, we have $\mathbf{m}(\mathbf{x})_{k^*} \geq \mathbf{m}(\mathbf{x})_k + 2\epsilon$. Finally, for any $k \neq k^*$, we get

$$\mathbf{m}(\mathbf{x} + \boldsymbol{\tau})_{k^*} \geq \mathbf{m}(\mathbf{x})_{k^*} - \epsilon \geq \mathbf{m}(\mathbf{x})_k + \epsilon \geq \mathbf{m}(\mathbf{x} + \boldsymbol{\tau})_k.$$

Then, $\operatorname{argmax}_{k \in [K]} \mathbf{m}(\mathbf{x})_k = \operatorname{argmax}_{k \in [K]} \mathbf{m}(\mathbf{x} + \boldsymbol{\tau})_k$. This concludes the proof. \square

Similarly, we can demonstrate a mode preservation property for robust classifiers w.r.t. the Renyi divergence.

Proposition 11 (Mode preservation for Renyi-robust classifiers) *Let $\mathbf{m} \in \mathcal{M}_\beta(\alpha_p, \epsilon)$ be a robust randomized classifier and $\mathbf{x} \in \mathcal{X}$ such that*

$$(\mathbf{m}(\mathbf{x})_{(1)})^{\frac{\beta}{\beta-1}} \geq \exp\left(2 - \frac{1}{\beta}\right)\epsilon \left(\mathbf{m}(\mathbf{x})_{(2)}\right)^{\frac{\beta-1}{\beta}}.$$

Then, for any $\boldsymbol{\tau} \in \mathcal{X}$, the following holds,

$$\|\boldsymbol{\tau}\|_p \leq \alpha_p \implies \mathbf{h}_{\text{rob}}(\mathbf{x}) = \mathbf{h}_{\text{rob}}(\mathbf{x} + \boldsymbol{\tau}),$$

where $\mathbf{h}_{\text{rob}}(\mathbf{x}) := \operatorname{argmax}_{k \in [K]} \mathbf{m}(\mathbf{x})_k$.

Proof Let $\mathbf{x}, \boldsymbol{\tau} \in \mathcal{X}$ such that $\|\boldsymbol{\tau}\|_p \leq \alpha_p$ and $m \in \mathcal{M}_\beta(\alpha_p, \epsilon)$ such that

$$(m(\mathbf{x})_{(1)})^{\frac{\beta}{\beta-1}} \geq \exp\left(\left(2 - \frac{1}{\beta}\right)\epsilon\right) (m(\mathbf{x})_{(2)})^{\frac{\beta-1}{\beta}}.$$

Then by definition of $\mathcal{M}_\beta(\alpha_p, \epsilon)$, we have

$$D_\beta(m(\mathbf{x}), m(\mathbf{x} + \boldsymbol{\tau})) \leq \epsilon.$$

Furthermore, by using Proposition 1, for any $k \in \{1, \dots, K\}$ we have

$$(*)m(\mathbf{x})_k \leq (\exp(\epsilon)m(\mathbf{x} + \boldsymbol{\tau})_k)^{\frac{\beta-1}{\beta}} \text{ and } (**)m(\mathbf{x} + \boldsymbol{\tau})_k \leq (\exp(\epsilon)m(\mathbf{x})_k)^{\frac{\beta-1}{\beta}}.$$

Let us denote k^* the index such that $m(\mathbf{x})_{k^*} = m(\mathbf{x})_{(1)}$. Then using (*) we get

$$m(\mathbf{x} + \boldsymbol{\tau})_{k^*} \geq \exp(-\epsilon)(m(\mathbf{x})_{k^*})^{\frac{\beta}{\beta-1}}.$$

Furthermore for any $k \in \{1, \dots, K\}$ where $k \neq k^*$, we can use the assumption we made on m to get

$$\exp(-\epsilon)(m(\mathbf{x})_{k^*})^{\frac{\beta}{\beta-1}} \geq \exp\left(\frac{\beta-1}{\beta}\epsilon\right) (m(\mathbf{x})_k)^{\frac{\beta-1}{\beta}}.$$

Finally, using (**) we have

$$\exp\left(\frac{\beta-1}{\beta}\epsilon\right) (m(\mathbf{x})_k)^{\frac{\beta-1}{\beta}} \geq m(\mathbf{x} + \boldsymbol{\tau})_k.$$

The above gives us $\operatorname{argmax}_{k \in [K]} m(\mathbf{x})_k = \operatorname{argmax}_{k \in [K]} m(\mathbf{x} + \boldsymbol{\tau})_k$. This concludes the proof. □

Coming back to the decomposition in Eq. (5), with the above result, we can bound the risk the adversary induces with non-zero perturbations by the mass of points on which the classifier \mathbf{h}_{rob} gives the good response but based on a low probability of success, i.e. with small confidence

$$\mathcal{R}_{>0}^{\text{adv}}(m) \leq \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}}[\mathbf{h}_{\text{rob}}(\mathbf{x}) = y \text{ and } m(\mathbf{x})_{(1)} < m(\mathbf{x})_{(2)} + 2\epsilon]. \tag{18}$$

This means that the only points on which the adversary may induce misclassification are the points on which m already has a high risk. Once more, this says something fundamental about the behavior of robust randomized classifiers. On undefended models, the adversary could change the decision on any point it wanted; now it is limited to changing points on which the classifier is already inaccurate. This considerably mitigates the threat model we should consider. Furthermore, for any deterministic classifier designed as in Eq. (17), we can also bound the maximal loss of accuracy under attack the classifier may suffer. This bound may, however, be harder to evaluate since it now depends on both the classifier and the dataset distribution. The classifier we define in Eq. (17) and the mode preservation property of m are closely related to provable defenses based on randomized smoothing. The core idea of randomized smoothing is to take a hypothesis \mathbf{h} and to build a robust classifier that writes

$$c_{rob} : \mathbf{x} \mapsto \operatorname{argmax}_{k \in [K]} \mathbb{P}_{z \sim \mathcal{N}(0, \sigma^2 I)} [\mathbf{h}(\mathbf{x} + \mathbf{z}) = k]. \quad (19)$$

From a probabilistic point of view, for any input \mathbf{x} , randomized smoothing amounts to output the most probable class of the probability measure $m(\mathbf{x}) := \mathbf{h} \# \mathcal{N}(\mathbf{x}, \sigma^2 I)$. Hence, randomized smoothing uses the mode preservation property of m to build a provably robust (deterministic) classifier. Therefore, the above results (Proposition 10 and Eq. 18) also hold for provable defenses based on randomized smoothing. Studying randomized smoothing from our point of view could give an interesting new perspective on that method. So far no results have been published on the generalisation gap of this defense in the adversarial setting. We could devise generalization bounds by similarity with our analysis. Furthermore, the probabilistic interpretation stresses that randomized smoothing is somewhat restrictive since it only considers probability measures which are the expectation on a simple noise injection scheme. The mode preservation property explains the behavior of randomized smoothing, but also presents fundamental properties of randomized defenses that could be used to construct more general defense schemes.

8 Numerical validations against ℓ_2 adversary

To illustrate our findings, we train randomized neural networks with Gaussian pre-processing during training and inference on CIFAR-10 and CIFAR-100. Based on this randomized classifier, we study the impact of randomization on the standard accuracy of the network, and observe the theoretical trade-off between accuracy and robustness.

8.1 Architecture and training procedure

All the neural networks we use in this section are WideResNets (Zagoruyko & Komodakis, 2016) with 28 layers, a widen factor of 10, a dropout factor of 0.3 and LeakyRelu activation with a 0.1 slope. To train an undefended standard classifier we use the following hyper-parameters.²

- *Number of Epochs:* 200
- *Batch size:* 400
- *Loss function:* Cross Entropy Loss
- *Optimizer:* Stochastic gradient descent algorithm with momentum 0.9, weight decay of 2×10^{-4} and a learning rate that decreases during the training as follows:

$$lr = \begin{cases} 0.1 & \text{if } 0 \leq \text{epoch} < 60 \\ 0.02 & \text{if } 60 \leq \text{epoch} < 120 \\ 0.004 & \text{if } 120 \leq \text{epoch} < 160 \\ 0.0008 & \text{if } 160 \leq \text{epoch} < 200. \end{cases}$$

To transform these standard networks into randomized classifiers, we inject noise drawn from Gaussian distributions, each with various standard deviations directly on the image

² Reusable code can be found in the following repository: <https://github.com/MILES-PSL/Adversarial-Robustness-Through-Randomization>.

Fig. 1 Impact of the standard deviation of the Gaussian noise on accuracy in a randomized model on CIFAR-10 and CIFAR-100 dataset

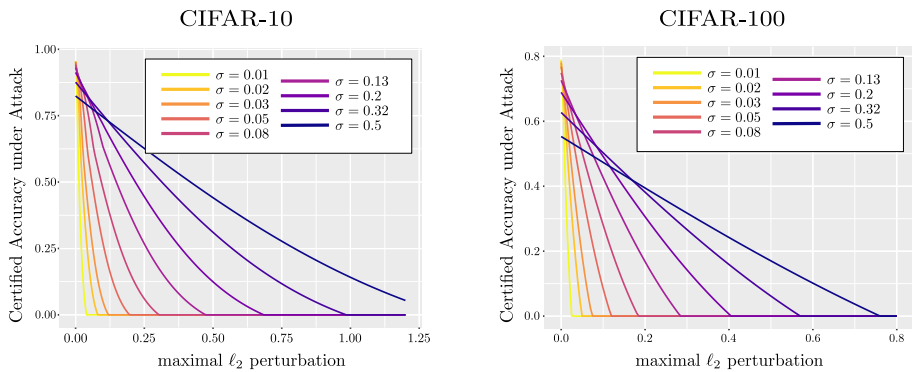
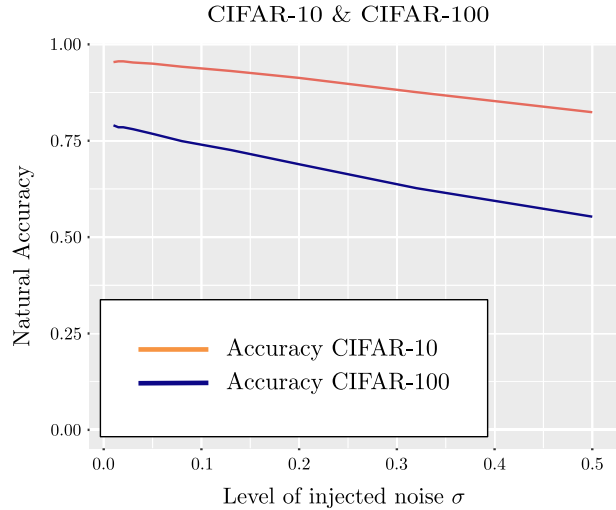


Fig. 2 Guaranteed accuracy of different randomized models with Gaussian noise given the ℓ_2 norm of the adversarial perturbations

before passing it through the network. Both during training and test, for computational efficiency, we evaluate the performance of the algorithm over a single run for every images; hence no Monte Carlo estimator is used. However, in practice, the test-time accuracy is stable when evaluated over the entire test dataset.

8.2 Results

Figures 1 and 2 show the accuracy and the minimum level of accuracy under attack of our randomized neural network for several levels of injected noise. We can see (Fig. 1) that the precision decreases as the noise intensity grows. In that sense, the noise must be calibrated to preserve both accuracy and robustness against adversarial attacks. This is to be expected, because the greater the entropy of the classifier, the less precise it gets.

Furthermore, when injecting Gaussian noise as a defense mechanism, the resulting randomized network m is both $(\alpha_2, \frac{(\alpha_2)^2}{2\sigma})$ -robust w.r.t. D_1 and $(\alpha_2, 2\Phi(\frac{\alpha_2}{2\sigma}) - 1)$ -robust w.r.t. D_{TV} against ℓ_2 adversaries. Therefore thanks to Theorems 3 and 5 we have that

$$\mathcal{R}^{\text{adv}}(m; \alpha_2) - \mathcal{R}(m) \leq 2\Phi\left(\frac{\alpha_2}{2\sigma}\right) - 1, \text{ and} \quad (20)$$

$$\mathcal{R}^{\text{adv}}(m; \alpha_2) - \mathcal{R}(m) \leq 1 - e^{-\frac{(\alpha_2)^2}{2\sigma}} \mathbb{E}_{x \sim \mathcal{D}_1} [e^{-H(m(x))}]. \quad (21)$$

Figure 2 illustrates the theoretical lower bound on accuracy under attack [based on the minimum gap between Eqs. (20) and (21)] for different standard deviations. The term in entropy has been estimated using a Monte Carlo method with 10^4 simulations. The trade-off between accuracy and robustness appears with respect to the noise intensity. With small noises, the accuracy is high, but the guaranteed accuracy drops fast with respect to the magnitude of the adversarial perturbation. Conversely, with bigger noises, the accuracy is lower but decreases slowly with respect to the magnitude of the adversarial perturbation. Overall, we get strong accuracy guarantees against small adversarial perturbations, but when the perturbation is bigger than 0.5 on CIFAR-10 (resp. 0.3 on CIFAR-100, the guarantees are still not sufficient).

9 Lesson learned and future work

This paper brings new contributions to the theory of robustness to adversarial attacks. We provided an in depth analysis of randomized classifier, demonstrating their interest to defend against adversarial attacks. We first defined a notion of robustness for randomized classifiers using probability metrics/divergences, namely the total variation distance and the Renyi divergence. Second, we demonstrated that when a randomized classifier complies with this definition of robustness, we can bound their loss of accuracy under attack. We also studied the generalization properties of this class of functions and gave results indicating that robust randomized classifiers can generalize. Finally, we showed that randomized classifiers have a mode preservation property. This presents a fundamental property of randomized defenses that can be used to explain randomized smoothing from a probabilistic point of view. To support our theoretical findings we presented a simple yet efficient scheme for building robust randomized classifiers. We show that Gaussian noise injection can provide principled robustness against ℓ_2 adversarial attacks. We ran a set of experiments on CIFAR-10 and CIFAR-100 using Gaussian noise injection with advanced neural network architectures to build accurate models with controlled loss of accuracy under attack.

Future work will focus on studying the combination of randomization with more sophisticated defenses and on devising new tight bounds on the adversarial generalization and the adversarial risk gap of randomized classifiers. Based on the connections we established we randomized smoothing in Sect. 7, we will also aim at devising bounds on the gap between the standard and adversarial risks for this defense. Another interesting direction would be to show that the classifiers based on randomized smoothing have a generalization gap similar to the classes of randomized classifiers we studied.

Appendix 1: Proof of technical lemmas

Appendix 1.1: Proof of Lemma 1

Proof Let $\beta > 1$. Let us denote g and g' respectively the probability density functions of ρ and ρ' with respect to the Lebesgue measure. We also set $\mathbf{x}' = \mathbf{x} + \boldsymbol{\tau}$ for readability. Then we have

$$\begin{aligned} D_\beta(\rho, \rho') &= \frac{1}{\beta - 1} \log \mathbb{E}_{\mathbf{z} \sim \rho'} \left[\left(\frac{g(\mathbf{z})}{g'(\mathbf{z})} \right)^\beta \right] \\ &= \frac{1}{\beta - 1} \log \mathbb{E}_{\mathbf{z} \sim \rho'} \left[\exp \left(\frac{\beta}{2} \left((\mathbf{z} - \mathbf{x}')^\top \Sigma^{-1} (\mathbf{z} - \mathbf{x}') - (\mathbf{z} - \mathbf{x})^\top \Sigma^{-1} (\mathbf{z} - \mathbf{x}) \right) \right) \right]. \end{aligned}$$

By change of variable we get

$$\begin{aligned} &= \frac{1}{\beta - 1} \log \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(0, \Sigma)} \left[\exp \left(\frac{\beta}{2} \left(\mathbf{z}^\top \Sigma^{-1} \mathbf{z} - (\mathbf{z} + \boldsymbol{\tau})^\top \Sigma^{-1} (\mathbf{z} + \boldsymbol{\tau}) \right) \right) \right] \\ &= \frac{1}{\beta - 1} \log \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(0, \Sigma)} \left[\exp \left(\frac{\beta}{2} \left(-2\mathbf{z}^\top \Sigma^{-1} \boldsymbol{\tau} - \|\boldsymbol{\tau}\|_{\Sigma^{-1}}^2 \right) \right) \right] \\ &= \frac{1}{\beta - 1} \log \int_{\mathbb{R}^d} \frac{\exp \left(-\frac{1}{2} \mathbf{z}^\top \Sigma^{-1} \mathbf{z} - \frac{\beta}{2} 2\mathbf{z}^\top \Sigma^{-1} \boldsymbol{\tau} - \frac{\beta}{2} \|\boldsymbol{\tau}\|_{\Sigma^{-1}}^2 \right)}{(2\pi)^d \det(\Sigma)^{d/2}} d\mathbf{z}. \end{aligned}$$

Furthermore, for any $\mathbf{z} \in \mathbb{R}^d$, we have

$$\begin{aligned} &-\frac{1}{2} \mathbf{z}^\top \Sigma^{-1} \mathbf{z} - \frac{\beta}{2} 2\mathbf{z}^\top \Sigma^{-1} \boldsymbol{\tau} - \frac{\beta}{2} \|\boldsymbol{\tau}\|_{\Sigma^{-1}}^2 \\ &= -\frac{1}{2} (\mathbf{z} + \beta\boldsymbol{\tau})^\top \Sigma^{-1} (\mathbf{z} + \beta\boldsymbol{\tau}) + \frac{\beta^2 - \beta}{2} \|\boldsymbol{\tau}\|_{\Sigma^{-1}}^2. \end{aligned}$$

Then we can rewrite the Renyi divergence as follows

$$\begin{aligned} D_\beta(\rho, \rho') &= \frac{1}{\beta - 1} \log \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(-\beta\boldsymbol{\tau}, \Sigma)} \left[\exp \left(\frac{\beta^2 - \beta}{2} \|\boldsymbol{\tau}\|_{\Sigma^{-1}}^2 \right) \right] \\ &= \frac{1}{\beta - 1} \log \left(\exp \left(\frac{\beta^2 - \beta}{2} \|\boldsymbol{\tau}\|_{\Sigma^{-1}}^2 \right) \right) \\ &= \frac{\beta}{2} \|\boldsymbol{\tau}\|_{\Sigma^{-1}}^2. \end{aligned}$$

This concludes the proof. □

Appendix 1.2: Proof of Lemma 2

Proof Let us denote g and g' respectively the probability density functions of ρ and ρ' with respect to the Lebesgue measure. Furthermore, we denote $\mathbf{x}' = \mathbf{x} + \boldsymbol{\tau}$. Then by definition of

the total variation distance, we have $D_{TV}(\rho, \rho') = \rho(Z) - \rho'(Z)$ with $Z = \{z \text{ s.t. } g(z) \geq g'(z)\}$. In our case $g(z) \geq g'(z)$ is equivalent to

$$(z - x')^T \Sigma^{-1} (z - x') - (z - x)^T \Sigma^{-1} (z - x) \geq 0.$$

Then with the same simplification as above, we have

$$\begin{aligned} \rho(Z) &= \mathbb{P}_{z \sim \mathcal{N}(x, \Sigma)} \left((z - x')^T \Sigma^{-1} (z - x') - (z - x)^T \Sigma^{-1} (z - x) \geq 0 \right) \\ &= \mathbb{P}_{z \sim \mathcal{N}(0, \Sigma)} \left((z - \tau)^T \Sigma^{-1} (z - \tau) - z^T \Sigma^{-1} z \geq 0 \right) \\ &= \mathbb{P}_{z \sim \mathcal{N}(0, \Sigma)} \left(-2z^T \Sigma^{-1} \tau + \|\tau\|_{\Sigma^{-1}}^2 \geq 0 \right) \\ &= \mathbb{P}_{z \sim \mathcal{N}(0, I_d)} \left(z^T \Sigma^{-1/2} \tau \leq \frac{1}{2} \|\tau\|_{\Sigma^{-1}}^2 \right). \end{aligned}$$

Furthermore, if $z \sim \mathcal{N}(0, I_d)$ then $z^T \Sigma^{-1/2} \tau \sim \mathcal{N}(0, \|\tau\|_{\Sigma^{-1}}^2)$; hence we also have $\frac{z^T \Sigma^{-1/2} \tau}{\|\tau\|_{\Sigma^{-1}}} \sim \mathcal{N}(0, 1)$. Accordingly we get

$$\rho(Z) = \mathbb{P}_{z \sim \mathcal{N}(0, 1)} \left(z \leq \frac{1}{2} \|\tau\|_{\Sigma^{-1}} \right) = \Phi \left(\frac{1}{2} \|\tau\|_{\Sigma^{-1}} \right).$$

By symmetry we get that $\rho'(A) = 1 - \rho(A) = 1 - \Phi \left(\frac{1}{2} \|\tau\|_{\Sigma^{-1}} \right)$. We then get

$$D_{TV}(\mu, \nu) = 2\Phi \left(\frac{\|\tau\|_{\Sigma^{-1}}}{2} \right) - 1$$

which concludes the proof. □

Appendix 2: Discussion on probability metrics

As mentioned earlier in this paper, the choice of the metric/divergence is crucial as it characterizes the notion of adversarial robustness we are examining. We focus on the total variation distance and Renyi divergence, but the question of whether these metrics/divergences are more appropriate than others remains open. It should be noted, however, that our definition of robustness is monotonous depending on the metric/divergence we use.

Proposition 12 (Monotonicity of the robustness) *Let m be a randomized classifier, and let D and D' be two divergences/metrics on $\mathcal{P}(\mathcal{Y})$. If there exists a non decreasing function $f : \mathbb{R} \mapsto \mathbb{R}$ such that $\forall \rho, \rho' \in \mathcal{P}(\mathcal{Y}), D(\rho, \rho') \leq f(D'(\rho, \rho'))$, then the following assertion holds.*

$$m \text{ is } (\alpha_p, \epsilon)\text{-robust w.r.t. } D' \implies m \text{ is } (\alpha_p, f(\epsilon))\text{-robust w.r.t. } D.$$

The proof straightforwardly comes from the definition of robustness.

Proof Let us consider m a randomized classifier (α_p, ϵ) -robust w.r.t. D' . Then for any $x \sim \mathcal{D}$, and τ s.t. $\|\tau\|_p \leq \alpha_p$, since f is non decreasing, we have

$$D(m(x), m(x + \tau)) \leq f(D'(m(x), m(x + \tau))) \leq f(\epsilon).$$

Then m is $(\alpha_p, f(\epsilon))$ -robust w.r.t. D which concludes the proof. □

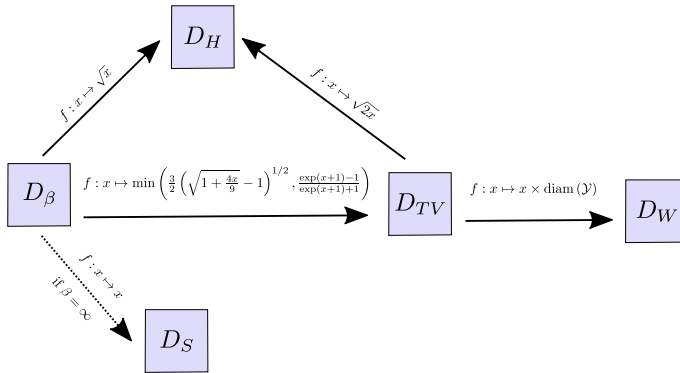


Fig. 3 Summary of the relations between the different robustness notions from Propositions 13 and 14

The above result suggests that the different notions of robustness we might conceive are more related than they appear. Here are some of the most classical divergences used in machine learning. Let ρ, ρ', ν three measures in $\mathcal{P}(\mathcal{Y})$. We denote g and g' the probability density functions of ρ and ρ' with respect to ν . Then we can define the *Wasserstein distance* as follows

$$D_W(\rho, \rho') := \inf \int_{\mathcal{Y}^2} \text{dist}(y, y') d\pi(y, y'), \tag{22}$$

where dist is some ground distance on \mathcal{Y} , and the infimum is taken over all joint distributions π in $\mathcal{P}(\mathcal{Y} \times \mathcal{Y})$ with marginals ρ and ρ' .

Remark 5 In transportation theory, the Wasserstein distance is solution of the Monge-Kantorovich problem with the cost function $c(y, y') = \text{dist}(y, y')$. Then, the definitions of total variation and Wasserstein distance match when we use the trivial distance $\text{dist}(y, y') = \mathbb{1}\{y \neq y'\}$.

We also define respectively the *Hellinger distance* and the *Separation distance* as follows.

$$D_H(\rho, \rho') := \left[\int_{\mathcal{Y}} (\sqrt{g} - \sqrt{g'})^2 d\nu \right]^{1/2}. \tag{23}$$

$$D_S(\rho, \rho') := \sup_{y \in \mathcal{Y}} \left(1 - \frac{g(y)}{g'(y)} \right). \tag{24}$$

If we take any of the above metrics/divergences to instantiate a notion of adversarial robustness we might get very different semantics for them. However, we can show that any of these definitions can be covered—with respect to Proposition 12—either by the Renyi or the total variation robustness. Figure 3 summarizes the links we can make between all these different definitions of robustness, and Propositions 13 and 14 present the associated

results. We can see that the total variation distance and the Renyi divergence are both central since they can cover any of the other robustness notions. This does not mean that they are more appropriate than the others, but at least they are general enough to cover a wide range of possible definitions.

Proposition 13 *Let m be a randomized classifier. If m is (α_p, ϵ) -robust w.r.t. D_{TV} then the following assertions hold.*

- m is $(\alpha_p, \epsilon \times \text{diam}(\mathcal{Y}))$ -robust w.r.t. D_W , where $\text{diam}(\mathcal{Y}) := \max_{y, y' \in \mathcal{Y}} \text{dist}(y, y')$.
- m is $(\alpha_p, \sqrt{2\epsilon})$ -robust w.r.t. D_H .

Proof Let us consider ρ and $\rho' \in \mathcal{P}(\mathcal{Y})$. Thanks to Gibbs and Su (2002) we have

- $D_W(\rho, \rho') \leq \text{diam}(\mathcal{Y})D_{TV}(\rho, \rho')$.
- $D_H(\rho, \rho') \leq \sqrt{2D_{TV}(\rho, \rho')}$.

Hence, by using Proposition 12 respectively with $f : x \mapsto \text{diam}(\mathcal{Y})x$ and $f : x \mapsto \sqrt{2x}$ we get the expected results. □

Proposition 14 *Let m be a randomized classifier. If m is (α_p, ϵ) -robust w.r.t. D_β then the following assertions hold.*

- m is (α_p, ϵ') -robust w.r.t. D_{TV} with $\epsilon' = \min\left(\frac{3}{2}\left(\sqrt{1 + \frac{4\epsilon}{9}} - 1\right)^{1/2}, \frac{\exp(\epsilon+1)-1}{\exp(\epsilon+1)+1}\right)$.
- m is $(\alpha_p, \sqrt{\epsilon})$ -robust w.r.t. D_H .
- If $\beta = \infty$, then m is (α_p, ϵ) robust w.r.t. D_S .

Proof (1) First, let us suppose that $\beta \geq 1$. Thanks to Proposition 2 and to (Gibbs & Su, 2002), for any $\rho, \rho' \in \mathcal{P}(\mathcal{Y})$ we have

- $D_H(\rho, \rho') \leq \sqrt{D_1(\rho, \rho')} \leq \sqrt{D_\beta(\rho, \rho')}$ (see Gibbs & Su, 2002).
- $D_{TV}(\rho, \rho') \leq \min\left(\frac{3}{2}\left(\sqrt{1 + \frac{4D_\beta(\rho, \rho')}{9}} - 1\right)^{1/2}, \frac{\exp(D_\beta(\rho, \rho')+1)-1}{\exp(D_\beta(\rho, \rho')+1)+1}\right)$ (Proposition 2).

Hence, by using Proposition 12, as above, we get the expected results.

(2) Now let us suppose that $\beta = \infty$. By definition of the supremum divergence, we have

$$D_\infty(\rho, \rho') = \sup_{B \subset \mathcal{Y}} \left| \ln \frac{\rho(B)}{\rho'(B)} \right|.$$

Furthermore, note that the function $x \mapsto 1 - x - |\ln(x)|$ is negative on \mathbb{R} , therefore for any $y \in \mathcal{Y}$ one has

$$1 - \frac{\rho(y)}{\rho'(y)} \leq \left| \ln \frac{\rho(y)}{\rho'(y)} \right|.$$

Since the above inequality is true for any $y \in \mathcal{Y}$, we have

$$D_S(\rho, \rho') = \sup_{y \in \mathcal{Y}} \left(1 - \frac{\rho(y)}{\rho'(y)} \right) \leq \sup_{y \in \mathcal{Y}} \left| \ln \frac{\rho(y)}{\rho'(y)} \right| \leq \sup_{B \subset \mathcal{Y}} \left| \ln \frac{\rho(B)}{\rho'(B)} \right| = D_\infty(\rho, \rho').$$

Finally, by using Proposition 12 with $f : x \mapsto x$ we get the expected results. \square

Author contributions All authors contributed to the problem definition. The technical content and proofs were mainly carried out by RP and LM under the joint supervision of FY, CG-P, JA and YC. The first version of the manuscript was written by RP and all authors commented on the previous versions of the manuscript. All authors read and approved the final manuscript.

Funding Open access funding provided by EPFL Lausanne. Rafael Pinot has been supported in part by Ecocloud, an EPFL research center (Postdoctoral Research Award). Finally, this work was granted access to OpenPOWER prototype from GENCI-IDRIS under the Preparatory Access AP010610510, and HPC resources of IDRIS under the allocation 2020-101141 made by GENCI.

Data availability We only use public (benchmark) dataset and our code is accessible online.

Code availability For direct access to the implementation, one can refer to the following Github repository <https://github.com/MILES-PSL/Adversarial-Robustness-Through-Randomization>.

Declarations

Conflict of interest The authors have no relevant financial or non-financial interests to disclose, no conflicts of interest to declare that are relevant to the content of this article.

Ethics approval All authors certify that they have no affiliations with or involvement in any organization or entity with any financial interest or non-financial interest in the subject matter or materials discussed in this manuscript nor financial or proprietary interests in any material discussed in this article. The authors approve the ethical standards of the publisher.

Consent to participate and publication All authors are aware of the submission of this manuscript and agree to its publication.

Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

References

- Athalye, A., Carlini, N., & Wagner, D. (2018). Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. In J. Dy & A. Krause (Eds.), *Proceedings of the 35th international conference on machine learning, proceedings of machine learning research* (vol. 8, pp. 274–283). Stockholm Sweden: Stockholm mässan.
- Awasthi, P., Frank, N., & Mohri, M. (2020). Adversarial learning guarantees for linear hypotheses and neural networks. In H. D. III & A. Singh (Eds.), *Proceedings of the 37th international conference on machine learning, proceedings of machine learning research* (vol. 119, pp. 431–441). PMLR.

- Bartlett, P. L., & Mendelson, S. (2002). Rademacher and Gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3, 463–482.
- Ben-Tal, A., El Ghaoui, L., & Nemirovski, A. (2009). *Robust optimization* (Vol. 28). Princeton University Press.
- Biggio, B., Corona, I., Maiorca, D., Nelson, B., Šrđnić, N., Laskov, P., Giacinto, G., & Roli, F. (2013). Evasion attacks against machine learning at test time. In *Joint European conference on machine learning and knowledge discovery in databases* (pp. 387–402). Springer.
- Bu, Z., Dong, J., Long, Q., & Weijie, S. (2020). Deep learning with Gaussian differential privacy. *Harvard Data Science Review*. <https://doi.org/10.1162/99608f92.cfc5dd25>.
- Carlini, N., & Wagner, D. (2017). Adversarial examples are not easily detected: Bypassing ten detection methods. In *Proceedings of the 10th ACM workshop on artificial intelligence and security* (pp. 3–14).
- Chapeau-Blondeau, F., & Rousseau, D. (2004). Noise-enhanced performance for an optimal Bayesian estimator. *IEEE Transactions on Signal Processing*, 52(5), 1327–1334. <https://doi.org/10.1109/TSP.2004.826176>.
- Chen, P. Y., Sharma, Y., Zhang, H., Yi, J., & Hsieh, C. J. (2018). Ead: Elastic-net attacks to deep neural networks via adversarial examples. In *AAAI*.
- Cohen, J., Rosenfeld, E., & Kolter, Z. (2019). Certified adversarial robustness via randomized smoothing. In K. Chaudhuri & R. Salakhutdinov (Eds.), *Proceedings of the 36th international conference on machine learning, proceedings of machine learning research* (vol. 97, pp. 1310–1320). PMLR.
- Cover, T. M., & Thomas, J. A. (2012). *Elements of information theory*. Wiley.
- Croce, F., & Hein, M. (2020). Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks. In *Proceedings of the 37th international conference on machine learning, ICML 2020, 13–18 July 2020, virtual event, proceedings of machine learning research* (vol. 119, pp. 2206–2216). PMLR.
- Dalvi, N., Domingos, P., Sanghai, S., & Verma, D. (2004). Adversarial classification. In *Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining* (pp. 99–108).
- Dhillon, G. S., Azizzadenesheli, K., Lipton, Z. C., Bernstein, J., Kossaifi, J., Khanna, A., & Anandkumar, A. (2018). Stochastic activation pruning for robust adversarial defense. In *6th international conference on learning representations, ICLR 2018, Vancouver, BC, Canada, April 30–May 3, 2018, conference track proceedings*. OpenReview.net.
- Gibbs, A. L., & Su, F. E. (2002). On choosing and bounding probability metrics. *International Statistical Review/Revue Internationale de Statistique*, 70(3), 419–435.
- Gilardi, G. L. (2010). On Pinsker's and Vajda's type inequalities for Csiszár's f -divergences. *IEEE Transactions on Information Theory*, 56(11), 5377–5386. <https://doi.org/10.1109/TIT.2010.2068710>.
- Globerson, A., & Roweis, S. (2006). Nightmare at test time: Robust learning by feature deletion. In *Proceedings of the 23rd international conference on Machine learning* (pp. 353–360).
- Goodfellow, I. J., Shlens, J., & Szegedy, C. (2015). Explaining and harnessing adversarial examples. In Y. Bengio & Y. LeCun (Eds.), *3rd International conference on learning representations, ICLR 2015, San Diego, CA, USA, May 7–9, 2015, conference track proceedings*.
- Grandvalet, Y., Canu, S., & Boucheron, S. (1997). Noise injection: Theoretical prospects. *Neural Computation*, 9(5), 1093–1108.
- He, W., Wei, J., Chen, X., Carlini, N., & Song, D. (2017). Adversarial example defense: Ensembles of weak defenses are not strong. In *11th USENIX Workshop on Offensive Technologies (WOOT 17)*.
- He, Z., Rakin, A. S., & Fan, D. (2019). Parametric noise injection: Trainable randomness to improve deep neural network robustness against adversarial attack. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 588–597).
- Hu, S., Yu, T., Guo, C., Chao, W. L., & Weinberger, K. Q. (2019). A new defense against adversarial images: Turning a weakness into a strength. In *Advances in neural information processing systems* (pp. 1635–1646).
- Jetley, S., Lord, N. A., & Torr, P. H. (2018). With friends like these, who needs adversaries? In *Proceedings of the 32nd international conference on neural information processing systems, NIPS'18, Red Hook, NY, USA* (pp. 10772–10782). Curran Associates Inc.
- Kearns, M., & Li, M. (1993). Learning in the presence of malicious errors. *SIAM Journal on Computing*, 22(4), 807–837.
- Kearns, M. J., Schapire, R. E., & Sellie, L. M. (1994). Toward efficient agnostic learning. *Machine Learning*, 17(2–3), 115–141.
- Khim, J., & Loh, P. L. (2018). Adversarial risk bounds for binary classification via function transformation. arXiv preprint [arXiv:1810.09519](https://arxiv.org/abs/1810.09519).

- Krizhevsky, A., & Hinton, G. (2009). *Learning multiple layers of features from tiny images*. Citeseer: Technical report.
- Langlois, A., Stehlé, D., & Steinfeld, R. (2014). Gghlite: More efficient multilinear maps from ideal lattices. In *Annual international conference on the theory and applications of cryptographic techniques* (pp. 239–256). Springer.
- Lecuyer, M., Atlidakis, V., Geambasu, R., Hsu, D., & Jana, S. (2019). Certified robustness to adversarial examples with differential privacy. In *2019 IEEE symposium on security and privacy (SP)* (pp. 656–672). IEEE.
- Li, B., Chen, C., Wang, W., & Carin, L. (2019). Certified adversarial robustness with additive noise. In *Advances in neural information processing systems* (pp. 9464–9474).
- Liu, X., Cheng, M., Zhang, H., & Hsieh, C. J. (2018). Towards robust neural networks via random self-ensemble. In *European conference on computer vision* (pp. 381–397). Springer.
- Lowd, D., & Meek, C. (2005). Adversarial learning. In *Proceedings of the eleventh ACM SIGKDD international conference on Knowledge discovery in data mining* (pp. 641–647).
- Madry, A., Makelov, A., Schmidt, L., Tsipras, D., & Vladu, A. (2018). Towards deep learning models resistant to adversarial attacks. In *6th International conference on learning representations, ICLR 2018, Vancouver, BC, Canada, April 30–May 3, 2018, conference track proceedings*. OpenReview.net.
- Metzen, J. H., Genewein, T., Fischer, V. & Bischoff, B. (2017). On detecting adversarial perturbations. In *5th International conference on learning representations, ICLR 2017, Toulon, France, April 24–26, 2017, conference track proceedings*. OpenReview.net.
- Mitaim, S., & Kosko, B. (1998). Adaptive stochastic resonance. *Proceedings of the IEEE*, 86(11), 2152–2183.
- Mohri, M., Rostamizadeh, A., & Talwalkar, A. (2018). *Foundations of machine learning*.
- Papernot, N., McDaniel, P., Wu, X., Jha, S., & Swami, A. (2016). Distillation as a defense to adversarial perturbations against deep neural networks. In *2016 IEEE symposium on security and privacy (SP)* (pp. 582–597). IEEE.
- Perez, L., & Wang, J. (2017). The effectiveness of data augmentation in image classification using deep learning. arXiv preprint [arXiv:1712.04621](https://arxiv.org/abs/1712.04621).
- Petzka, H., Kamp, M., Adilova, L., Sminchisescu, C., & Boley, M. (2021). Relative flatness and generalization. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P. Liang, & J. W. Vaughan (Eds.), *Advances in neural information processing systems* (Vol. 34, pp. 18420–18432). Curran Associates Inc.
- Peyré, G., & Cuturi, M. (2019). Computational optimal transport: With applications to data science. *Foundations and Trends® in Machine Learning* 11(5-6): 355–607.
- Pinot, R., Meunier, L., Araujo, A., Kashima, H., Yger, F., Gouy-Pailler, C., & Atif, J. (2019). Theoretical evidence for adversarial robustness through randomization. In *Advances in neural information processing systems* (pp. 11838–11848).
- Rényi, A. (1961). *On measures of entropy and information*. Hungarian Academy of Sciences Budapest Hungary: Technical report.
- Robert, C. (2007). *The Bayesian choice: From decision-theoretic foundations to computational implementation*. Springer.
- Salman, H., Li, J., Razenshteyn, I., Zhang, P., Zhang, H., Bubeck, S., & Yang, G. (2019). Provably robust deep learning via adversarially trained smoothed classifiers. In *Advances in neural information processing systems* (pp. 11289–11300).
- Schmidt, L., Santurkar, S., Tsipras, D., Talwar, K., & Madry, A. (2018). Adversarially robust generalization requires more data. In *Advances in neural information processing systems* (pp. 5014–5026).
- Shalev-Shwartz, S., & Ben-David, S. (2014). *Understanding machine learning: From theory to algorithms*. Cambridge University Press.
- Sharif, M., Bhagavatula, S., Bauer, L., & Reiter, M. K. (2016). Accessorize to a crime: Real and stealthy attacks on state-of-the-art face recognition. In *Proceedings of the 2016 ACM SIGSAC conference on computer and communications security* (pp. 1528–1540).
- Simon-Gabriel, C. J., Ollivier, Y., Bottou, L., Schölkopf, B., & Lopez-Paz, D. (2019). First-order adversarial vulnerability of neural networks and input dimension. In *International conference on machine learning* (pp. 5809–5817).
- Sitawarin, C., Bhagoji, A. N., Mosenia, A., Chiang, M., & Mittal, P. (2018). Darts: Deceiving autonomous cars with toxic signs. arXiv preprint [arXiv:1802.06430](https://arxiv.org/abs/1802.06430).
- Su, D., Zhang, H., Chen, H., Yi, J., Chen, P. Y., & Gao, Y. (2018). Is robustness the cost of accuracy?—A comprehensive study on the robustness of 18 deep image classification models. In *Proceedings of the European conference on computer vision (ECCV)* (pp. 631–648).

- Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I. J., & Fergus, R. (2014). Intriguing properties of neural networks. In Y. Bengio & Y. LeCun (Eds.), *2nd international conference on learning representations, ICLR 2014, Banff, AB, Canada, April 14–16, 2014, conference track proceedings*.
- Tramer, F., Carlini, N., Brendel, W., & Madry, A. (2020). On adaptive attacks to adversarial example defenses. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, & H. Lin (Eds.), *Advances in neural information processing systems* (Vol. 33, pp. 1633–1645). Curran Associates Inc.
- Tsipras, D., Santurkar, S., Engstrom, L., Turner, A., & Madry, A. (2019). Robustness may be at odds with accuracy. In *7th International conference on learning representations, ICLR 2019, New Orleans, LA, USA, May 6–9, 2019*. OpenReview.net.
- Vajda, I. (1970). Note on discrimination information and variation. *IEEE Transactions on Information Theory*, 16(6), 771–773.
- Van der Vaart, A. W. (2000). *Asymptotic statistics* (Vol. 3). Cambridge University Press.
- van Erven, T., & Harremoës, P. (2014). Rényi divergence and Kullback–Leibler divergence. *IEEE Transactions on Information Theory*, 60(7), 3797–3820.
- Verma, G., & Swami, A. (2019). Error correcting output codes improve probability estimation and adversarial robustness of deep neural networks. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, & R. Garnett (Eds.), *Advances in neural information processing systems 32* (pp. 8646–8656). Curran Associates Inc.
- Villani, C. (2003). *Topics in optimal transportation*. Number 58. American Mathematical Soc.
- Xie, C., Wang, J., Zhang, Z., Ren, Z., & Yuille, A. L. (2018). Mitigating adversarial effects through randomization. In *6th International conference on learning representations, ICLR 2018, Vancouver, BC, Canada, April 30–May 3, 2018, conference track proceedings*. OpenReview.net.
- Xu, H., & Mannor, S. (2012). Robustness and generalization. *Machine Learning*, 86(3), 391–423.
- Yang, G., Duan, T., Hu, E., Salman, H., Razenshteyn, I., & Li, J. (2020). Randomized smoothing of all shapes and sizes.
- Yao, D., Xi, Z., Tianyi, Z., Chen, C., Guannan, L., & Miryung, K. (2020). An analysis of adversarial attacks and defenses on autonomous driving models. In *18th Annual IEEE international conference on pervasive computing and communications*. IEEE.
- Yin, D., Kannan, R., & Bartlett, P. (2019). Rademacher complexity for adversarially robust generalization. In *International conference on machine learning* (pp. 7085–7094).
- Zagoruyko, S., & Komodakis, N. (2016). Wide residual networks. In *Proceedings of the British machine vision conference (BMVC)* (pp. 87.1–87.12). BMVA Press.
- Zhang, H., Yu, Y., Jiao, J., Xing, E. P., Ghaoui, L. E., & Jordan, M. I. (2019). Theoretically principled trade-off between robustness and accuracy. In K. Chaudhuri & R. Salakhutdinov (Eds.), *Proceedings of the 36th international conference on machine learning, ICML 2019, 9–15 June 2019, Long Beach, California, USA, proceedings of machine learning research* (vol. 97, pp. 7472–7482). PMLR.
- Zozor, S., & Amblard, P. O. (1999). Stochastic resonance in discrete time nonlinear AR(1) models. *IEEE Transactions on Signal Processing*, 47(1), 108–122.