



**HAL**  
open science

# A Kernel Perspective for Regularizing Deep Neural Networks

Alberto Bietti, Grégoire Mialon, Dexiong Chen, Julien Mairal

► **To cite this version:**

Alberto Bietti, Grégoire Mialon, Dexiong Chen, Julien Mairal. A Kernel Perspective for Regularizing Deep Neural Networks. ICML 2019 - 36th International Conference on Machine Learning, Jun 2019, Long Beach, United States. pp.664-674. hal-01884632v4

**HAL Id: hal-01884632**

**<https://inria.hal.science/hal-01884632v4>**

Submitted on 14 May 2019

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

---

# A Kernel Perspective for Regularizing Deep Neural Networks

---

Alberto Bietti<sup>\*1</sup> Grégoire Mialon<sup>\*12</sup> Dexiong Chen<sup>1</sup> Julien Mairal<sup>1</sup>

## Abstract

We propose a new point of view for regularizing deep neural networks by using the norm of a reproducing kernel Hilbert space (RKHS). Even though this norm cannot be computed, it admits upper and lower approximations leading to various practical strategies. Specifically, this perspective (i) provides a common umbrella for many existing regularization principles, including spectral norm and gradient penalties, or adversarial training, (ii) leads to new effective regularization penalties, and (iii) suggests hybrid strategies combining lower and upper bounds to get better approximations of the RKHS norm. We experimentally show this approach to be effective when learning on small datasets, or to obtain adversarially robust models.

## 1. Introduction

Learning predictive models for complex tasks often requires large amounts of annotated data. For instance, convolutional neural networks are huge-dimensional and typically involve more parameters than training samples, which raises several challenges: achieving good generalization with small datasets is indeed difficult, which limits the deployment of such deep models to many tasks where labeled data is scarce, *e.g.*, in biology (Ching et al., 2018). Besides, imperceptible adversarial perturbations can significantly degrade the prediction quality (Szegedy et al., 2013; Biggio & Roli, 2018). These issues raise the question of regularization as an essential tool to control the complexity of deep models, as well as their stability to small variations of their inputs.

In this paper, we present a new perspective on regularization of deep networks, by viewing convolutional neural networks (CNNs) as elements of a RKHS following the work of Bietti & Mairal (2019) on deep convolutional kernels. For

---

<sup>\*</sup>Equal contribution <sup>1</sup>Univ. Grenoble Alpes, Inria, CNRS, Grenoble INP, LJK, 38000 Grenoble, France <sup>2</sup>Département d’informatique de l’ENS, ENS, CNRS, Inria, PSL, 75005 Paris, France. Correspondence to: <firstname.lastname@inria.fr>.

such kernels, the RKHS contains indeed deep convolutional networks similar to generic ones—up to smooth approximations of rectified linear units. Such a point of view provides a natural regularization function, the RKHS norm, which allows us to control the variations of the predictive model and to limit its complexity for better generalization. Besides, the norm also acts as a Lipschitz constant, which provides a direct control on the stability to adversarial perturbations.

In contrast to traditional kernel methods, the RKHS norm cannot be explicitly computed in our setup. Yet, this norm admits numerous approximations—lower bounds and upper bounds—which lead to many strategies for regularization based on penalties, constraints, or combinations thereof. Depending on the chosen approximation, we recover then many existing principles such as spectral norm regularization (Cisse et al., 2017; Yoshida & Miyato, 2017; Miyato et al., 2018a; Sedghi et al., 2019), gradient penalties and double backpropagation (Drucker & Le Cun, 1991; Simon-Gabriel et al., 2019; Gulrajani et al., 2017; Roth et al., 2017; 2018; Arbel et al., 2018), adversarial training (Madry et al., 2018), and we also draw links with tangent propagation (Simard et al., 1998). For all these principles, we provide a unified viewpoint and theoretical insights, and we also introduce new variants, which we show are effective in practice when learning with few labeled data, or in the presence of adversarial perturbations.

Moreover, regularization and robustness are tightly linked in our kernel framework. Specifically, some lower bounds on the RKHS norm lead to robust optimization objectives with worst-case  $\ell_2$  perturbations; further, we can extend margin-based generalization bounds in the spirit of Bartlett et al. (2017); Boucheron et al. (2005) to the setting of *adversarially robust* generalization (see Schmidt et al., 2018), where an adversary can perturb test data. We also discuss connections between recent regularization strategies for training generative adversarial networks and approaches to generative modeling based on kernel two-sample tests (MMD) (Dzugaite et al., 2015; Li et al., 2017; Bińkowski et al., 2018).

## Summary of the contributions.

- We introduce an RKHS perspective for regularizing deep neural networks models which provides a unified view on various practical regularization principles, together with theoretical insight and guarantees;

- By considering lower bounds to the RKHS norm, we obtain new penalties based on adversarial perturbations, adversarial deformations, or gradient norms of prediction functions, which we show to be effective in practice;
- Our RKHS point of view suggests combined strategies based on both upper and lower bounds, which we show often perform empirically best in the context of generalization from small image and biological datasets, by providing a tighter control of the RKHS norm.

**Related work.** The construction of hierarchical kernels and the study of neural networks in the corresponding RKHS was studied by Mairal (2016); Zhang et al. (2016; 2017); Bietti & Mairal (2019). Some of the regularization strategies we obtain from our kernel perspective are variants of previous approaches to adversarial robustness (Cisse et al., 2017; Madry et al., 2018; Simon-Gabriel et al., 2019; Roth et al., 2018), to improving generalization (Drucker & Le Cun, 1991; Miyato et al., 2018b; Sedghi et al., 2019; Simard et al., 1998; Yoshida & Miyato, 2017), and stable training of generative adversarial networks (Roth et al., 2017; Gulrajani et al., 2017; Arbel et al., 2018; Miyato et al., 2018a). The link between robust optimization and regularization was studied by Xu et al. (2009a;b), focusing mainly on linear models with quadratic or hinge losses. The notion of adversarial generalization was considered by Schmidt et al. (2018), who provide lower bounds on a particular data distribution. Sinha et al. (2018) provide generalization guarantees in the different setting of distributional robustness; compared to our bound, they consider expected loss instead of classification error, and their bounds do not highlight the dependence on the model complexity.

## 2. Regularization of Deep Neural Networks

In this section, we recall the kernel perspective on deep networks introduced by Bietti & Mairal (2019), and present upper and lower bounds on the RKHS norm of a given model, leading to various regularization strategies. For simplicity, we first consider real-valued networks and binary classification, before discussing multi-class extensions.

### 2.1. Relation between deep networks and RKHSs

Kernel methods consist of mapping data living in a set  $\mathcal{X}$  to a RKHS  $\mathcal{H}$  associated to a positive definite kernel  $K$  through a mapping function  $\Phi : \mathcal{X} \rightarrow \mathcal{H}$ , and then learning simple machine learning models in  $\mathcal{H}$ . Specifically, when considering a real-valued regression or binary classification problem, classical kernel methods find a prediction function  $f : \mathcal{X} \rightarrow \mathbb{R}$  living in the RKHS which can be written in linear form, i.e., such that  $f(x) = \langle f, \Phi(x) \rangle_{\mathcal{H}}$  for all  $x$  in  $\mathcal{X}$ . While explicit mapping to a possibly infinite-dimensional space is of course only an abstract mathematical operation, learning  $f$  can be done implicitly by com-

puting kernel evaluations and typically by using convex programming (Schölkopf & Smola, 2001).

Moreover, the RKHS norm  $\|f\|_{\mathcal{H}}$  acts as a natural regularization function, which controls the variations of model predictions according to the geometry induced by  $\Phi$ :

$$|f(x) - f(x')| \leq \|f\|_{\mathcal{H}} \cdot \|\Phi(x) - \Phi(x')\|_{\mathcal{H}}. \quad (1)$$

Unfortunately, our setup does not allow us to use the RKHS norm in a traditional way since evaluating the kernel is intractable. Instead, we propose a different approach that considers explicit parameterized representations of functions contained in the RKHS, given by generic CNNs, and leverage properties of the RKHS and the kernel mapping in order to regularize when learning the network parameters.

Consider indeed a real-valued deep convolutional network  $f : \mathcal{X} \rightarrow \mathbb{R}$ , where  $\mathcal{X}$  is simply  $\mathbb{R}^d$ , with rectified linear unit (ReLU) activations and no bias units. By constructing an appropriate multi-layer hierarchical kernel, Bietti & Mairal (2019) show that the corresponding RKHS  $\mathcal{H}$  contains a CNN with the same architecture and parameters as  $f$ , but with activations that are smooth approximations of ReLU. Although the model predictions might not be strictly equal, we will abuse notation and denote this approximation with smooth ReLU by  $f$  as well, with the hope that the regularization procedures derived from the RKHS model will be effective in practice on the original CNN  $f$ .

Besides, the mapping  $\Phi(\cdot)$  is shown to be non-expansive:

$$\|\Phi(x) - \Phi(x')\|_{\mathcal{H}} \leq \|x - x'\|_2, \quad (2)$$

so that controlling  $\|f\|_{\mathcal{H}}$  provides some robustness to additive  $\ell_2$ -perturbations, by (1). Additionally, with appropriate pooling operations, Bietti & Mairal (2019) show that the kernel mapping is also stable to deformations, meaning that the RKHS norm also controls robustness to translations and other transformations including scaling and rotations, which can be seen as deformations when they are small.

In contrast to standard kernel methods, where the RKHS norm is typically available in closed form, this norm is difficult to compute in our setup, and requires approximations. The following sections present upper and lower bounds on  $\|f\|_{\mathcal{H}}$ , with linear convolutional operations denoted by  $W_k$  for  $k = 1, \dots, L$ , where  $L$  is the number of layers. Defining  $\theta := \{W_k : k = 1, \dots, L\}$ , we then leverage these bounds to approximately solve the following penalized or constrained optimization problems on a training set  $(x_i, y_i), i = 1, \dots, n$ :

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^n \ell(y_i, f_{\theta}(x_i)) + \lambda \|f_{\theta}\|_{\mathcal{H}}^2 \quad \text{or} \quad (3)$$

$$\min_{\theta: \|f_{\theta}\|_{\mathcal{H}} \leq C} \frac{1}{n} \sum_{i=1}^n \ell(y_i, f_{\theta}(x_i)). \quad (4)$$

We also note that while the construction of [Bietti & Mairal \(2019\)](#) considers VGG-like networks ([Simonyan & Zisserman, 2014](#)), the regularization algorithms we obtain in practice can be easily adapted to different architectures such as residual networks ([He et al., 2016](#)).

## 2.2. Exploiting lower bounds of the RKHS norm

In this section, we devise regularization algorithms by leveraging lower bounds on  $\|f\|_{\mathcal{H}}$ , obtained by relying on the following variational characterization of Hilbert norms:

$$\|f\|_{\mathcal{H}} = \sup_{\|u\|_{\mathcal{H}} \leq 1} \langle f, u \rangle_{\mathcal{H}}.$$

At first sight, this definition is not useful since the set  $U = \{u \in \mathcal{H} : \|u\|_{\mathcal{H}} \leq 1\}$  may be infinite-dimensional and the inner products  $\langle f, u \rangle_{\mathcal{H}}$  cannot be computed in general. Thus, we devise tractable lower bound approximations by considering smaller sets  $\bar{U} \subset U$ .

**Adversarial perturbation penalty.** Thanks to the non-expansiveness of  $\Phi$ , we can consider the subset  $\bar{U} \subset U$  defined as  $\bar{U} = \{\Phi(x + \delta) - \Phi(x) : x \in \mathcal{X}, \|\delta\|_2 \leq 1\}$ , leading to the bound

$$\|f\|_{\mathcal{H}} \geq \|f\|_{\delta}^2 := \sup_{x \in \mathcal{X}, \|\delta\|_2 \leq 1} f(x + \delta) - f(x), \quad (5)$$

which is reminiscent of adversarial perturbations. Adding a regularization parameter  $\epsilon > 0$  in front of the norm then corresponds to different sizes of perturbations:

$$\epsilon \|f\|_{\mathcal{H}} = \sup_{\|u\|_{\mathcal{H}} \leq \epsilon} \langle f, u \rangle_{\mathcal{H}} \geq \sup_{x \in \mathcal{X}, \|\delta\|_2 \leq \epsilon} f(x + \delta) - f(x). \quad (6)$$

Using this lower bound or its square as a penalty in the objective (3) when training a CNN provides a way to regularize. Optimizing over adversarial perturbations has been useful to obtain robust models (e.g., the PGD method of [Madry et al., 2018](#)); yet our approach differs in two important ways:

(i) it involves a penalty that is decoupled from the loss term such that in principle, our penalty could be used beyond the supervised empirical risk paradigm. In contrast, PGD optimizes the robust formulation (7) below, which fits training data while considering perturbations on the loss.

(ii) our penalty involves a global maximization problem on the input space  $\mathcal{X}$ , as opposed to only maximizing on perturbations near training data. In practice, optimizing over  $\mathcal{X}$  is however difficult and instead, we replace  $\mathcal{X}$  by random mini-batches of examples, yielding further lower bounds on the RKHS norm. These examples may be labeled or not, in contrast to PGD that perturb labeled examples only. When using such a mini-batch, a gradient of the penalty can be obtained by first finding maximizers  $\hat{x}, \hat{\delta}$  (where  $\hat{x}$  is an element of the mini-batch and  $\hat{\delta}$  is a perturbation), and then

computing gradients of  $f_{\theta}(\hat{x} + \hat{\delta}) - f_{\theta}(\hat{x})$  with respect to  $\theta$  by using back-propagation. In practice, we compute the perturbations  $\delta$  for each example  $x$  by using a few steps of projected gradient ascent with constant step-lengths.

**Robust optimization yields another lower bound.** In some contexts, our penalized approach is related to solving the robust optimization problem

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^n \sup_{\|\delta\|_2 \leq \epsilon} \ell(y_i, f_{\theta}(x_i + \delta)), \quad (7)$$

which is commonly considered for training adversarially robust classifiers ([Wong & Kolter, 2018](#); [Madry et al., 2018](#); [Raghunathan et al., 2018](#)). In particular, [Xu et al. \(2009b\)](#) show that the penalized and robust objectives are equivalent in the case of the hinge loss with linear predictors, when the data is non-separable. They also show the equivalence for kernel methods when considering the (intractable) full perturbation set  $U$  around each point in the RKHS  $\Phi(x_i)$ , that is, predictions  $\langle f, \Phi(x_i) + u \rangle_{\mathcal{H}}$  with  $u \in U$ . Intuitively, when a training example  $(x_i, y_i)$  is misclassified, we are in the “linear” part of the hinge loss, such that

$$\sup_{\|u\|_{\mathcal{H}} \leq \epsilon} \ell(y_i, \langle f, \Phi(x_i) + u \rangle_{\mathcal{H}}) = \ell(y_i, f(x_i)) + \epsilon \|f\|_{\mathcal{H}}.$$

For other losses such as the logistic loss, a regularization effect is still present even for correctly classified examples, though it may be smaller since the loss has a reduced slope for such points. This leads to an *adaptive* regularization mechanism that may automatically reduce the amount of regularization when the data is easily separable. However, the robust optimization approach might only encourage local stability around training examples, while the global quantity  $\|f\|_{\mathcal{H}}$  may become large in order to better fit the data. We note that a perfect fit of the data with large complexity does not prevent generalization (see, e.g., [Belkin et al., 2018a;b](#)); yet, such mechanisms are still poorly understood. Nevertheless, it is easy to show that the robust objective (7) lower bounds the penalized objective with penalty  $\epsilon \|f\|_{\mathcal{H}}$ .

**Gradient penalties.** Taking  $\bar{U} = \{\frac{\Phi(x) - \Phi(y)}{\|x - y\|_2} : x, y \in \mathcal{X}\}$ , which is a subset of  $U$  by Eq. (2)—it turns out that this is the same set as for adversarial perturbation penalties, since  $\Phi$  is homogeneous ([Bietti & Mairal, 2019](#)) and  $\mathcal{X} = \mathbb{R}^d$ —we obtain a lower bound based on the Lipschitz constant of  $f$ :

$$\|f\|_{\mathcal{H}} \geq \sup_{x, y \in \mathcal{X}} \frac{f(x) - f(y)}{\|x - y\|_2} \geq \|\nabla f\| := \sup_{x \in \mathcal{X}} \|\nabla f(x)\|_2, \quad (8)$$

where the second inequality becomes an equality when  $\mathcal{X}$  is convex, and the supremum is taken over points where  $f$  is differentiable. Although we are unaware of previous work using this exact lower bound for a generic regularization



penalty, we note that variants replacing the supremum over  $x$  by an expectation over data have been recently used to stabilize the training of generative adversarial networks (Gulrajani et al., 2017; Roth et al., 2017), and we provide insights in Section 3.2 on the benefits of RKHS regularization in such a setting. Related penalties have been considered in the context of robust optimization, for regularization or robustness, noting that a penalty based on the gradient of the loss function  $x \mapsto \ell(y, f(x))$  can give a good approximation of (7) when  $\epsilon$  is small (Drucker & Le Cun, 1991; Lyu et al., 2015; Roth et al., 2018; Simon-Gabriel et al., 2019).

**Penalties based on deformation stability.** We may also obtain new penalties by considering more exotic sets  $\bar{U} = \{\Phi(\tilde{x}) - \Phi(x) : x \in \mathcal{X}, \tilde{x} \text{ is a small deformation of } x\}$ , where the amount of deformation is dictated by the stability bounds of Bietti & Mairal (2019) in order to ensure that  $\bar{U} \subset U$ . More precisely, such bounds depend on the maximum displacement and Jacobian norm of the diffeomorphisms considered. These can be easily computed for various parameterized families of transformations, such as translations, scaling or rotations, leading to simple ways to control the regularization strength through the parameters of these transformations. One can also consider infinitesimal deformations from such parameterized transformations, which approximately yields the *tangent propagation* regularization strategy of Simard et al. (1998). These approaches are detailed in Appendix B. If instead we consider the robust optimization formulation (7), we obtain a form of *data augmentation* where transformations are optimized instead of sampled, as done by (Engstrom et al., 2017).

**Extensions to multiple classes and beyond** We now extend the regularization strategies based on lower bounds to multi-valued networks, in order to deal with multiple classes. For that purpose, we consider a multi-class penalty  $\|f_1\|_{\mathcal{H}}^2 + \dots + \|f_K\|_{\mathcal{H}}^2$  for an  $\mathbb{R}^K$ -valued function  $f = (f_1, f_2, \dots, f_K)$ , and we define

$$\|f\|_{\delta}^2 := \sum_{k=1}^K \|f_k\|_{\delta}^2 \quad \text{and} \quad \|\nabla f\|^2 := \sum_{k=1}^K \|\nabla f_k\|^2,$$

where  $\|f_k\|_{\delta}$  is the adversarial penalty (5), and  $\|\nabla f_k\|$  is defined in (8). For deformation stability penalties, we proceed in a similar manner, and for robust optimization formulations (7), the extension is straightforward, given that multi-class losses such as cross-entropy can be directly optimized in an adversarial training or gradient penalty setup.

Finally, we note that while the kernel approach we introduce considers the Euclidian geometry in the input space, it is possible to consider heuristic alternatives for other geometries, such as  $\ell_{\infty}$  perturbations, as discussed in Appendix D.

### 2.3. Exploiting upper bounds with spectral norms

Instead of lower bounds, one may use instead the following upper bound from Bietti & Mairal (2019, Proposition 14):

$$\|f\|_{\mathcal{H}} \leq \omega(\|W_1\|, \dots, \|W_L\|), \quad (9)$$

where  $\omega$  is increasing in all of its arguments, and  $\|W_k\|$  is the spectral norm of the linear operator  $W_k$ . Here, we simply consider the spectral norm on the filters, given by  $\|W\| := \sup_{\|x\|_2 \leq 1} \|Wx\|_2$ . Other generalization bounds relying on similar quantities have been proposed for controlling complexity (Bartlett et al., 2017; Neyshabur et al., 2018), suggesting that using them for regularization is relevant even beyond our kernel perspective, as observed by Cisse et al. (2017); Sedghi et al. (2019); Yoshida & Miyato (2017). Extensions to multiple classes are simple to obtain by simply considering spectral norms up to the last layer.

**Penalizing the spectral norms.** One way to control the upper bound (9) when learning a neural network  $f_{\theta}$  is to consider a regularization penalty based on spectral norms

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^n \ell(y_i, f_{\theta}(x_i)) + \lambda \sum_{l=1}^L \|W_l\|^2, \quad (10)$$

where  $\lambda$  is a regularization parameter. To optimize this cost, one can obtain (sub)gradients of the penalty by computing singular vectors associated to the largest singular value of each  $W_l$ . We consider the method of Yoshida & Miyato (2017), which computes such singular vectors approximately using one or two iterations of the power method, as well as a more costly approach using the full SVD.

**Constraining the spectral norms with a continuation approach.** In the constrained setting, we want to optimize:

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^n \ell(y_i, f_{\theta}(x_i)) \quad \text{s.t.} \quad \|W_l\| \leq \tau; l \in 1, \dots, L,$$

where  $\tau$  is a user-defined constraint. This objective may be optimized by projecting each  $W_l$  in the spectral norm ball of radius  $\tau$  after each gradient step. Such a projection is achieved by truncating the singular values to be smaller than  $\tau$  (see Appendix C). We found that the loss was hardly optimized with this approach, and therefore introduce a continuation approach with an exponentially decaying schedule for  $\tau$  reaching a constant  $\tau_0$  after a few epochs, which we found to be important for good empirical performance.

### 2.4. Combining upper and lower bounds.

One advantage of lower bound penalties is that they are independent of the model parameterization, making them flexible enough to use with more complex architectures. In

addition, the connection with robust optimization can provide a useful mechanism for adaptive regularization. However, they do not provide a guaranteed control on the RKHS norm, unlike the upper bound strategies. This is particularly true for robust optimization approaches, which may favor small training loss and local stability over global stability through  $\|f\|_{\mathcal{H}}$ . Nevertheless, we observed that our new approaches based on separate penalties sometimes do help in controlling upper bounds as well (see Section 4).

While these upper bound strategies are useful for limiting model complexity, we found them empirically less effective for robustness (see Section 4.2). However, we observed that combining with lower bound approaches can overcome this weakness, perhaps due to a better control of local stability. In particular, such combined approaches often provide the best generalization performance in small data scenarios, as well as better guarantees on adversarially robust generalization thanks to a tighter control of the RKHS norm.

### 3. Theoretical Guarantees and Insights

In this section, we study how the kernel perspective allows us to extend standard margin-based generalization bounds to an adversarial setting in order to provide theoretical guarantees on adversarially robust generalization. We then discuss how our kernel approach provides novel interpretations for training generative adversarial networks.

#### 3.1. Guarantees on adversarial generalization

While various methods have been introduced to empirically gain robustness to adversarial perturbations, the ability to generalize with such perturbations, also known as *adversarial generalization* (Schmidt et al., 2018), still lacks theoretical understanding. Margin-based bounds have been useful to explain the generalization behavior of learning algorithms that can fit the training data well, such as kernel methods, boosting and neural networks (Koltchinskii & Panchenko, 2002; Boucheron et al., 2005; Bartlett et al., 2017). Here, we show how such arguments can be adapted to obtain guarantees on adversarial generalization, *i.e.*, on the expected classification error in the presence of an  $\ell_2$ -bounded adversary, based on the RKHS norm of a learned model. For a binary classification task with labels in  $\mathcal{Y} = \{-1, 1\}$  and data distribution  $\mathcal{D}$ , we would like to bound the expected adversarial error of a classifier  $f$ , given for some  $\epsilon > 0$  by

$$\text{err}_{\mathcal{D}}(f, \epsilon) := P_{(x,y) \sim \mathcal{D}}(\exists \|\delta\|_2 \leq \epsilon : yf(x + \delta) < 0). \quad (11)$$

Leveraging the fact that  $f$  is  $\|f\|_{\mathcal{H}}$ -Lipschitz, we now show how to further bound this quantity using empirical margins, following the usual approach to obtaining margin bounds for kernel methods (*e.g.*, Boucheron et al., 2005). Consider a training dataset  $(x_1, y_1), \dots, (x_n, y_n) \in \mathcal{X} \times \mathcal{Y}$ . Defining

$L_n^\gamma(f) := \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{y_i f(x_i) < \gamma\}$ , we have the following bound, proved in Appendix E:

**Proposition 1** (Adversarially robust margin bound). *With probability  $1 - \delta$  over a dataset  $\{(x_i, y_i)\}_{i=1, \dots, n}$ , we have, for all choices of  $\gamma > 0$  and  $f \in \mathcal{H}$ ,*

$$\text{err}_{\mathcal{D}}(f, \epsilon) \leq L_n^{\gamma+2\epsilon\|f\|_{\mathcal{H}}}(f) + \tilde{O}\left(\frac{\|f\|_{\mathcal{H}} \bar{B}}{\gamma \sqrt{n}}\right), \quad (12)$$

where  $\bar{B} = \sqrt{\frac{1}{n} \sum_{i=1}^n K(x_i, x_i)}$  and  $\tilde{O}$  hides a term depending logarithmically on  $\|f\|_{\mathcal{H}}$ ,  $\gamma$ , and  $\delta$ .

When  $\epsilon = 0$ , we obtain the usual margin bound, while  $\epsilon > 0$  yields a bound on adversarial error  $\text{err}_{\mathcal{D}}(f, \epsilon)$ , for some neural network  $f$  learned from data. Note that other complexity measures based on products of spectral norms may be used instead of  $\|f\|_{\mathcal{H}}$ , as well as multi-class extensions, following Bartlett et al. (2017); Neyshabur et al. (2018). In concurrent work, Khim & Loh (2018); Yin et al. (2019) derive similar bounds in the context of fully-connected networks. In contrast to these works, which bound complexity of a modified function class, our bound uses the complexity of the original class and leverages smoothness properties of functions to derive the margin bound.

One can then study the effectiveness of a regularization algorithm by inspecting cumulative distribution (CDF) plots of the *normalized margins*  $\bar{\gamma}_i = y_i f(x_i) / \|f\|_{\mathcal{H}}$ , for different strengths of regularization (an example is given in Figure 2, Section 4.2). According to the bound (12), one can assess expected adversarial error with  $\epsilon$ -bounded perturbations by looking at the part of the plot to the right of  $\bar{\gamma} = 2\epsilon$ . In particular, the value of the CDF at such a value of  $\bar{\gamma}$  is representative of the bound for large  $n$  (since the second term is negligible), while for smaller  $n$ , the best bound is obtained for a larger value of  $\bar{\gamma}$ , which also suggests that the right side of the plots is indicative of performance on small datasets.

When the RKHS norm can be well approximated, our bound provides a certificate on test error in the presence of adversaries. While such an approximation is difficult to obtain in general, the guarantee is most useful when lower and upper bounds of the RKHS norm are controlled together.

#### 3.2. New insights on generative adversarial networks

Generative adversarial networks (GANs) attempt to learn a *generator* neural network  $G_\phi : \mathcal{Z} \rightarrow \mathcal{X}$ , so that the distribution of  $G_\phi(z)$  with  $z \sim D_z$  a noise vector resembles a data distribution  $D_x$ . In this section, we discuss connections between recent regularization techniques for training GANs, and approaches to learning generative models based on a MMD criterion (Gretton et al., 2012), in view of our RKHS framework. Our goal is to provide a new insight on these methods, but not necessarily to provide a new one.

Various recent approaches have relied on regularization strategies on a *discriminator* network in order to improve the stability of GAN training and the quality of the produced samples. Some of these resemble the approaches presented in Section 2 such as gradient penalties (Gulrajani et al., 2017; Roth et al., 2017) and spectral norm regularization (Miyato et al., 2018a). We provide an RKHS interpretation of these methods as optimizing an MMD distance with the convolutional kernel introduced in Section 2:

$$\min_{\phi} \sup_{\|f\|_{\mathcal{H}} \leq 1} \mathbb{E}_{x \sim D_x} [f(x)] - \mathbb{E}_{z \sim D_z} [f(G_{\phi}(z))]. \quad (13)$$

When learning from an empirical distribution over  $n$  samples, the MMD criterion is known to have much better sample complexity than the Wasserstein-1 distance considered by Arjovsky et al. (2017) for high-dimensional data such as images (Sriperumbudur et al., 2012). While the MMD approach has been used for training generative models, it generally relies on a generic kernel function, such as a Gaussian kernel, that appears explicitly in the objective (Dziugaite et al., 2015; Li et al., 2017; Bińkowski et al., 2018). Although using a learned feature extractor can improve this, the Gaussian kernel might be a poor choice when dealing with natural signals such as images, while the hierarchical kernel we consider in our paper is better suited for this type of data, by providing useful invariance and stability properties. Leveraging the variational form of the MMD (13) with this kernel suggests for instance using convolutional networks as the discriminator  $f$ , with constraints on the spectral norms in order to ensure  $\|f\|_{\mathcal{H}} \leq C$  for some  $C$ , as done by Miyato et al. (2018a) through normalization.

## 4. Experiments

We tested the regularization strategies presented in Section 2 in the context of improving generalization on small datasets and training robust models. Our goal is to use common architectures used for large datasets and improve their performance in different settings through regularization. Our Pytorch implementation of the various strategies is available at [https://github.com/albietz/kernel\\_reg](https://github.com/albietz/kernel_reg).

For the adversarial training strategies, the inner maximization problems are solved using 5 steps of projected gradient ascent with constant step-lengths. In the case of the lower bound penalties  $\|f\|_{\delta}^2$  and  $\|\nabla f\|^2$ , we also maximize over examples in the mini-batch, only considering the maximal element when computing gradients with respect to parameters. For the robust optimization problem (7), we use PGD with  $\ell_2$  perturbations, as well as the corresponding  $\ell_2$  (squared) gradient norm penalty on the loss. For the upper bound approaches with spectral norms (SNs), we consider the SN projection strategy with decaying  $\tau$ , as well as the SN penalty (10), either using power iteration (PI) or a full SVD for computing gradients.

Table 1. Regularization on CIFAR10 with 1 000 examples for VGG-11 and ResNet-18. Each entry shows the test accuracy with/without data augmentation when all hyper-parameters are optimized on a validation set. See also Section A.1 in the appendix for additional results and statistical testing.

| Method                       | 1k VGG-11                   | 1k ResNet-18                |
|------------------------------|-----------------------------|-----------------------------|
| No weight decay              | 50.70 / 43.75               | 45.23 / 37.12               |
| Weight decay                 | 51.32 / 43.95               | 44.85 / 37.09               |
| SN penalty (PI)              | 54.64 / 45.06               | 47.01 / 39.63               |
| SN projection                | 54.14 / <b>46.70</b>        | 47.12 / 37.28               |
| VAT                          | 50.88 / 43.36               | 47.47 / 42.82               |
| PGD- $\ell_2$                | 51.25 / 44.40               | 45.80 / 41.87               |
| grad- $\ell_2$               | <b>55.19</b> / 43.88        | <b>49.30</b> / <b>44.65</b> |
| $\ f\ _{\delta}^2$ penalty   | 51.41 / 45.07               | 48.73 / 43.72               |
| $\ \nabla f\ ^2$ penalty     | 54.80 / 46.37               | <b>48.99</b> / <b>44.97</b> |
| PGD- $\ell_2$ + SN proj      | 54.19 / <b>46.66</b>        | 47.47 / 41.25               |
| grad- $\ell_2$ + SN proj     | <b>55.32</b> / <b>46.88</b> | 48.73 / 42.78               |
| $\ f\ _{\delta}^2$ + SN proj | 54.02 / <b>46.72</b>        | 48.12 / 43.56               |
| $\ \nabla f\ ^2$ + SN proj   | <b>55.24</b> / <b>46.80</b> | <b>49.06</b> / <b>44.92</b> |

### 4.1. Improving generalization on small datasets

We consider the datasets CIFAR10 and MNIST when using a small number of training examples, as well as 102 datasets of biological sequences that suffer from small sample size.

**CIFAR10.** In this setting, we use 1 000 and 5 000 examples of the CIFAR10 dataset, with or without data augmentation. We consider a VGG network (Simonyan & Zisserman, 2014) with 11 layers, as well as a residual network (He et al., 2016) with 18 layers, which achieve 91% and 93% test accuracy respectively when trained on the full training set with standard data augmentation (horizontal flips + random crops). We do not use batch normalization layers in order to prevent any interaction with spectral norms. Each strategy derived in Section 2 is trained for 500 epochs using SGD with momentum and batch size 128, halving the step-size every 40 epochs. In order to study the potential effectiveness of each method, we assume that a reasonably large validation set is available to select hyper-parameters; thus, we keep 10 000 annotated examples for this purpose. We also show results using a smaller validation set in Appendix A.1.

Table 1 shows the test accuracies on 1 000 examples for upper and lower bound approaches, as well as combined ones. We also include virtual adversarial training (VAT, Miyato et al., 2018b). We provide extended tables in Appendix A.1 with additional methods, other geometries, results for 5 000 examples, as well as hypothesis tests for comparing pairs of methods and assessing the significance of our findings. Overall, we find that the combined lower bound + SN constraints approaches often yield better results than either method separately. For lower bound approaches alone, we found our  $\|f\|_{\delta}^2$  and  $\|\nabla f\|^2$  penalties to often work best, particularly without data augmentation, while robust optimization strategies can be preferable with data

Table 2. Regularization on 300 or 1 000 examples from MNIST, using deformations from Infinite MNIST. (\*) indicates that random deformations were included as training examples, while  $\|f\|_\tau^2$  and  $\|D_\tau f\|^2$  use them as part of the regularization penalty. See Section A.2 in the appendix for more results and statistical testing.

| Method   | 300 VGG      | 1k VGG       |
|--|--------------|--------------|
| Weight decay   | 89.32        | 94.08        |
| SN projection  | 90.69        | 95.01        |
| grad- $\ell_2$                                       | 93.63        | 96.67        |
| $\ f\ _\delta^2$ penalty                             | 94.17        | 96.99        |
| $\ \nabla f\ ^2$ penalty                             | 94.08        | 96.82        |
| Weight decay (*)                                     | 92.41        | 95.64        |
| grad- $\ell_2$ (*)                                   | 95.05        | 97.48        |
| $\ D_\tau f\ ^2$ penalty                             | 94.18        | 96.98        |
| $\ f\ _\tau^2$ penalty                               | 94.42        | 97.13        |
| $\ f\ _\tau^2 + \ \nabla f\ ^2$                      | 94.75        | 97.40        |
| $\ f\ _\tau^2 + \ f\ _\delta^2$                      | 95.23        | <b>97.66</b> |
| $\ f\ _\tau^2 + \ f\ _\delta^2$ (*)                  | <b>95.53</b> | <b>97.56</b> |
| $\ f\ _\tau^2 + \ f\ _\delta^2 + \text{SN proj}$     | 95.20        | <b>97.60</b> |
| $\ f\ _\tau^2 + \ f\ _\delta^2 + \text{SN proj}$ (*) | <b>95.40</b> | <b>97.77</b> |

augmentation, perhaps thanks to the adaptive regularization effect discussed earlier, which may be helpful in this easier setting. Gradient penalties often outperform adversarial perturbation strategies, possibly because of the closed form gradients which may improve optimization. We also found that adversarial training strategies tend to poorly control SNs compared to gradient penalties, particularly PGD (see also Section 4.2). SN constraints alone can also work well in some cases, particularly for VGG architectures, and often outperform SN penalties. SN penalties can work well nevertheless and provide computational benefits when using the power iteration variant.

**Infinite MNIST.** In order to assess the effectiveness of lower bound penalties based on deformation stability, we consider the Infinite MNIST dataset (Loosli et al., 2007), which provides an “infinite” number of transformed generated examples for each of the 60 000 MNIST training digits. Here, we use a 5-layer VGG-like network with average pooling after each 3x3 convolution layer, in order to more closely match the architecture assumptions of Bietti & Mairal (2019) for deformation stability. We consider two lower bound penalties that leverage the digit transformations in Infinite MNIST: one based on “adversarial” deformations around each digit, denoted  $\|f\|_\tau^2$ ; and a tangent propagation (Simard et al., 1998) variant, denoted  $\|D_\tau f\|^2$ , which provides an approximation to  $\|f\|_\tau^2$  for small deformations based on gradients along a few tangent vector directions given by deformations (see Appendix B for details). Table 2 shows the obtained test accuracy for subsets of MNIST of size 300 and 1 000. Overall, we find that combining both adversarial penalties  $\|f\|_\tau^2$  and  $\|f\|_\delta^2$  performs best, which suggests that it is helpful to obtain tighter lower approximations of the RKHS norm by considering perturbations of

Table 3. Regularization on protein homology detection tasks, with or without data augmentation (DA). Fixed hyperparameters are selected using the first half of the datasets, and we report the average auROC50 score on the second half. See Section A.3 in the appendix for more details and statistical testing.

| Method                     | No DA        | DA           |
|----------------------------|--------------|--------------|
| No weight decay            | 0.421        | 0.541        |
| Weight decay               | 0.432        | 0.544        |
| SN proj                    | 0.583        | 0.615        |
| PGD- $\ell_2$              | 0.488        | 0.554        |
| grad- $\ell_2$             | 0.551        | 0.570        |
| $\ f\ _\delta^2$           | 0.577        | 0.611        |
| $\ \nabla f\ ^2$           | 0.566        | 0.598        |
| PGD- $\ell_2$ + SN proj    | <b>0.615</b> | <b>0.622</b> |
| grad- $\ell_2$ + SN proj   | 0.581        | <b>0.634</b> |
| $\ f\ _\delta^2$ + SN proj | <b>0.631</b> | <b>0.639</b> |
| $\ \nabla f\ ^2$ + SN proj | 0.576        | <b>0.617</b> |

different kinds. Explicitly controlling the spectral norms can further improve performance, as does training on deformed digits, which may yield better margins by exploiting the additional knowledge that small deformations preserve labels. Note that data augmentation alone (with some weight decay) does quite poorly in this case, even compared to our lower bound penalties which do not use deformations.

**Protein homology detection.** Remote homology detection between protein sequences is an important problem to understand protein structure. Given a protein sequence, the goal is to predict whether it belongs to a superfamily of interest. We consider the Structural Classification Of Proteins (SCOP) version 1.67 dataset (Murzin et al., 1995), which we process as described in Appendix A.3 in order to obtain 102 balanced binary classification tasks with 100 protein sequences each, thus resulting in a low-sample regime. Protein sequences were also cut to 400 amino acids.

Sequences are represented with a one-hot encoding strategy—that is, a sequence of length  $l$  is represented as a binary matrix in  $\{0, 1\}^{20 \times l}$ , where 20 is the number of different amino acids (alphabet size of the sequences). Such a structure can then be processed by convolutional neural networks (Alipanahi et al., 2015). In this paper, we do not try to optimize the structure of the network for the task, since our goal is only to evaluate the effect of regularization strategies. Therefore, we use a simple convolutional network with 3 convolutional layers followed by global max-pooling and a final fully-connected layer (we use filters of size 5, and a max-pooling layer after the second convolutional layer).

Training was done using Adam with a learning rate fixed to 0.01, and a weight decay parameter tuned for each method. Since hyper-parameter selection per dataset is difficult due to the low sample size, we use the same parameters across datasets. This allows us to use the first 51 datasets as a validation set for hyper-parameter tuning, and we report



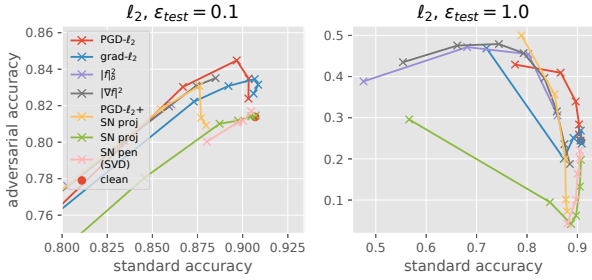


Figure 1. Robustness trade-off curves of different regularization methods for VGG11 on CIFAR10. Each plot shows test accuracy vs adversarial test accuracy for  $\ell_2$ -bounded, 40-step PGD adversaries with a fixed  $\epsilon_{test}$ . Different points on a curve correspond to training with different regularization strengths. The regularization increases monotonically along a given curve, and the leftmost points correspond to the strongest regularization. For PGD- $\ell_2$  + SN projection, we vary  $\epsilon$  with a fixed  $\tau = 0.8$ .

average performance with these fixed choices on the remaining 51 datasets. The standard performance measure for this task is the auROC50 score (area under the ROC curve up to 50% false positives). We note that the selection of hyperparameters has a transductive component, since some of the sequences in the test datasets may also appear in the datasets used for validation (possibly with a different label).

The results are shown in Table 3. The procedure used for data augmentation (right column) is described in Appendix A.3. We found that the most effective approach is the adversarial perturbation penalty, together with SN constraints. In particular, we found it to outperform the gradient penalty  $\|\nabla f\|^2$ , perhaps because in this case gradient penalties are only computed on a discrete set of possible points given by one-hot encodings, while adversarial perturbations may increase stability to wider regions, potentially covering different possible encoded sequences.

## 4.2. Training adversarially robust models

We consider the same VGG architecture as in Section 4.1, trained on CIFAR10 with data augmentation, with different regularization strategies. Each method is trained for 300 epochs using SGD with momentum and batch size 128, dividing the step-size in half every 30 epochs. This strategy was successful in reaching convergence for all methods.

Figure 1 shows the test accuracy of the different methods in the presence of  $\ell_2$ -bounded adversaries, plotted against standard accuracy. We can see that the robust optimization approaches tend to work better in high-accuracy regimes, perhaps because the local stability that they encourage is sufficient on this dataset, while the  $\|f\|_\delta^2$  penalty can be useful in large-perturbation regimes. We find that upper bound approaches alone do not provide robust models, but combining the SN constraint approach with a lower bound

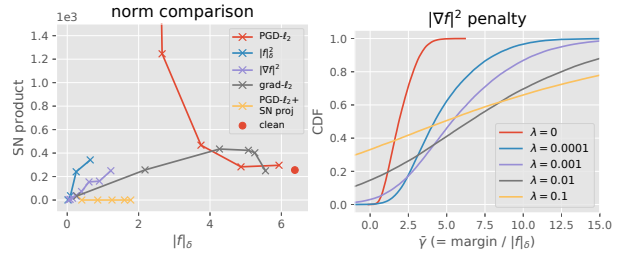


Figure 2. (left) Comparison of lower and upper bound quantities ( $\|f\|_\delta$  vs the product of spectral norms). (right) CDF plot of normalized empirical margins for the  $\|\nabla f\|^2$  penalty with different regularization strengths, normalized by  $\|f\|_\delta$ . We consider 1000 fixed training examples when computing  $\|f\|_\delta$ .

strategy (in this case PGD- $\ell_2$ ) helps improve robustness perhaps thanks to a more explicit control of stability. The plots also confirm that gradient penalties on the loss may be preferable for small regularization strengths (they achieve higher accuracy while improving robustness for small  $\epsilon_{test}$ ), while for stronger regularization, the gradient approximation no longer holds and the adversarial training approaches such as PGD (and its combination with SN constraints) are preferred. More experiments confirming these findings are available in Section A.4 of the appendix.

**Norm comparison and adversarial generalization.** Figure 2 (left) compares lower and upper bound quantities for different regularization strengths. Note that for PGD, in contrast to other methods, we can see that the product of spectral norms (representative of an upper bound on  $\|f\|_{\mathcal{H}}$ ) increases when the lower bound  $\|f\|_\delta$  decreases. This suggests that a network learned with PGD with large  $\epsilon$  may have large RKHS norm, possibly because the approach tries to separate  $\epsilon$ -balls around the training examples, which may require a more complex model than simply separating the training examples (see also Madry et al., 2018). This large discrepancy between upper and lower bounds highlights the fact that such models may only be stable locally near training data, though this happens to be enough for robustness on many test examples on CIFAR10.

In contrast, for other methods, and in particular the lower bound penalties  $\|f\|_\delta^2$  and  $\|\nabla f\|^2$ , the upper and lower bounds appear more tightly controlled, suggesting a more appropriate control of the RKHS norm. This makes our guarantees on adversarial generalization more meaningful, and thus we may look at the empirical distributions of normalized margins  $\bar{\gamma}$  obtained using  $\|f\|_\delta$  for normalization (as an approximation of  $\|f\|_{\mathcal{H}}$ ), shown in Figure 2 (right). The curves suggest that for small  $\bar{\gamma}$ , and hence small  $\epsilon_{test}$ , smaller values of  $\lambda$  are preferred, while stronger regularization helps for larger  $\bar{\gamma}$ , yielding lower test error guarantees in the presence of stronger adversaries according to our bounds in Section 3.1. This qualitative behavior is indeed observed in the results of Figure 1 on test data for the  $\|\nabla f\|^2$  penalty.

## Acknowledgements

This work was supported by the ERC grant number 714381 (SOLARIS project) and by the MSR-Inria joint centre.

## References

- Alipanahi, B., DeLong, A., Weirauch, M. T., and Frey, B. J. Predicting the sequence specificities of dna-and rna-binding proteins by deep learning. *Nature biotechnology*, 33(8):831, 2015.
- Altschul, S. F., Madden, T. L., Schäffer, A. A., Zhang, J., Zhang, Z., Miller, W., and Lipman, D. J. Gapped blast and psi-blast: a new generation of protein database search programs. *Nucleic acids research*, 25(17):3389–3402, 1997.
- Arbel, M., Sutherland, D. J., Bińkowski, M., and Gretton, A. On gradient regularizers for MMD GANs. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2018.
- Arjovsky, M., Chintala, S., and Bottou, L. Wasserstein GAN. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2017.
- Bartlett, P. L., Foster, D. J., and Telgarsky, M. J. Spectrally-normalized margin bounds for neural networks. In *Advances in Neural Information Processing Systems (NIPS)*, 2017.
- Belkin, M., Hsu, D., and Mitra, P. Overfitting or perfect fitting? risk bounds for classification and regression rules that interpolate. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2018a.
- Belkin, M., Ma, S., and Mandal, S. To understand deep learning we need to understand kernel learning. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2018b.
- Bietti, A. and Mairal, J. Group invariance, stability to deformations, and complexity of deep convolutional representations. *Journal of Machine Learning Research (JMLR)*, 20(25):1–49, 2019.
- Biggio, B. and Roli, F. Wild patterns: Ten years after the rise of adversarial machine learning. *Pattern Recognition*, 84:317–331, 2018.
- Bińkowski, M., Sutherland, D. J., Arbel, M., and Gretton, A. Demystifying MMD GANs. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2018.
- Boucheron, S., Bousquet, O., and Lugosi, G. Theory of classification: A survey of some recent advances. *ESAIM: probability and statistics*, 9:323–375, 2005.
- Ching, T. et al. Opportunities and obstacles for deep learning in biology and medicine. *Journal of The Royal Society Interface*, 15(141), 2018.
- Cisse, M., Bojanowski, P., Grave, E., Dauphin, Y., and Usunier, N. Parseval networks: Improving robustness to adversarial examples. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2017.
- Drucker, H. and Le Cun, Y. Double backpropagation increasing generalization performance. In *International Joint Conference on Neural Networks (IJCNN)*, 1991.
- Dziugaite, G. K., Roy, D. M., and Ghahramani, Z. Training generative neural networks via maximum mean discrepancy optimization. In *Conference on Uncertainty in Artificial Intelligence (UAI)*, 2015.
- Engstrom, L., Tsipras, D., Schmidt, L., and Madry, A. A rotation and a translation suffice: Fooling cnns with simple transformations. *arXiv preprint arXiv:1712.02779*, 2017.
- Gretton, A., Borgwardt, K. M., Rasch, M. J., Schölkopf, B., and Smola, A. A kernel two-sample test. *Journal of Machine Learning Research (JMLR)*, 13(Mar):723–773, 2012.
- Gulrajani, I., Ahmed, F., Arjovsky, M., Dumoulin, V., and Courville, A. C. Improved training of Wasserstein GANs. In *Advances in Neural Information Processing Systems (NIPS)*, 2017.
- Håndstad, T., Hestnes, A. J., and Sætrum, P. Motif kernel generated by genetic programming improves remote homology and fold detection. *BMC bioinformatics*, 8(1):23, 2007.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- Kakade, S. M., Sridharan, K., and Tewari, A. On the complexity of linear prediction: Risk bounds, margin bounds, and regularization. In *Advances in Neural Information Processing Systems (NIPS)*, 2009.
- Khim, J. and Loh, P.-L. Adversarial risk bounds via function transformation. *arXiv preprint arXiv:1810.09519*, 2018.
- Koltchinskii, V. and Panchenko, D. Empirical margin distributions and bounding the generalization error of combined classifiers. *The Annals of Statistics*, 30(1):1–50, 2002.
- Li, C.-L., Chang, W.-C., Cheng, Y., Yang, Y., and Póczos, B. Mmd gan: Towards deeper understanding of moment matching network. In *Advances in Neural Information Processing Systems (NIPS)*, 2017.

- Loosli, G., Canu, S., and Bottou, L. Training invariant support vector machines using selective sampling. In Bottou, L., Chapelle, O., DeCoste, D., and Weston, J. (eds.), *Large Scale Kernel Machines*, pp. 301–320. MIT Press, Cambridge, MA., 2007.
- Lyu, C., Huang, K., and Liang, H.-N. A unified gradient regularization family for adversarial examples. In *IEEE International Conference on Data Mining (ICDM)*, 2015.
- Madry, A., Makelov, A., Schmidt, L., Tsipras, D., and Vladu, A. Towards deep learning models resistant to adversarial attacks. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2018.
- Mairal, J. End-to-end kernel learning with supervised convolutional kernel networks. In *Advances in Neural Information Processing Systems (NIPS)*, 2016.
- Miyato, T., Kataoka, T., Koyama, M., and Yoshida, Y. Spectral normalization for generative adversarial networks. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2018a.
- Miyato, T., Maeda, S.-i., Ishii, S., and Koyama, M. Virtual adversarial training: a regularization method for supervised and semi-supervised learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI)*, 2018b.
- Murzin, A. G., Brenner, S. E., Hubbard, T., and Chothia, C. Scop: a structural classification of proteins database for the investigation of sequences and structures. *Journal of molecular biology*, 247(4):536–540, 1995.
- Neysshabur, B., Bhojanapalli, S., McAllester, D., and Srebro, N. A PAC-Bayesian approach to spectrally-normalized margin bounds for neural networks. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2018.
- Raghunathan, A., Steinhardt, J., and Liang, P. Certified defenses against adversarial examples. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2018.
- Roth, K., Lucchi, A., Nowozin, S., and Hofmann, T. Stabilizing training of generative adversarial networks through regularization. In *Advances in Neural Information Processing Systems (NIPS)*, 2017.
- Roth, K., Lucchi, A., Nowozin, S., and Hofmann, T. Adversarially robust training through structured gradient regularization. *arXiv preprint arXiv:1805.08736*, 2018.
- Schmidt, L., Santurkar, S., Tsipras, D., Talwar, K., and Madry, A. Adversarially robust generalization requires more data. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2018.
- Schölkopf, B. and Smola, A. J. *Learning with kernels: support vector machines, regularization, optimization, and beyond*. 2001.
- Sedghi, H., Gupta, V., and Long, P. M. The singular values of convolutional layers. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2019.
- Simard, P. Y., LeCun, Y. A., Denker, J. S., and Victorri, B. Transformation invariance in pattern recognition—tangent distance and tangent propagation. In *Neural networks: tricks of the trade*, pp. 239–274. Springer, 1998.
- Simon-Gabriel, C.-J., Ollivier, Y., Bottou, L., Schölkopf, B., and Lopez-Paz, D. First-order adversarial vulnerability of neural networks and input dimension. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2019.
- Simonyan, K. and Zisserman, A. Very deep convolutional networks for large-scale image recognition. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2014.
- Sinha, A., Namkoong, H., and Duchi, J. Certifying some distributional robustness with principled adversarial training. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2018.
- Sriperumbudur, B. K., Fukumizu, K., Gretton, A., Schölkopf, B., Lanckriet, G. R., et al. On the empirical estimation of integral probability metrics. *Electronic Journal of Statistics*, 6:1550–1599, 2012.
- Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., and Fergus, R. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013.
- Tsipras, D., Santurkar, S., Engstrom, L., Turner, A., and Madry, A. Robustness may be at odds with accuracy. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2019.
- Wong, E. and Kolter, J. Z. Provable defenses against adversarial examples via the convex outer adversarial polytope. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2018.
- Xu, H., Caramanis, C., and Mannor, S. Robust regression and lasso. In *Advances in Neural Information Processing Systems (NIPS)*, 2009a.
- Xu, H., Caramanis, C., and Mannor, S. Robustness and regularization of support vector machines. *Journal of Machine Learning Research (JMLR)*, 10(Jul):1485–1510, 2009b.

Yin, D., Ramchandran, K., and Bartlett, P. Rademacher complexity for adversarially robust generalization. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2019.

Yoshida, Y. and Miyato, T. Spectral norm regularization for improving the generalizability of deep learning. *arXiv preprint arXiv:1705.10941*, 2017.

Zhang, Y., Lee, J. D., and Jordan, M. I.  $\ell_1$ -regularized neural networks are improperly learnable in polynomial time. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2016.

Zhang, Y., Liang, P., and Wainwright, M. J. Convexified convolutional neural networks. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2017.



Section A of this supplementary presents extended results from our experiments, along with statistical tests for assessing the significance of our findings. Section B details our lower bound penalties based on deformations and their relationship to tangent propagation. Section C presents our continuation algorithm for optimization with spectral norm constraints. Section D describes heuristic extensions of our lower bound regularization strategies to non-Euclidian geometries. Finally, Section E provides our proof of the margin bound of Proposition 1 for adversarial generalization.

## A. Additional Experiment Results

### A.1. CIFAR10

This section provides more extensive results for the experiments on CIFAR10 from Section 4.1. In particular, Table 4 shows additional experiments on larger subsets of size 5000, as well as more methods, including different geometries (see Appendix D). The table also reports results obtained when using a smaller validation set of size 1000. The full hyper-parameter grid is given in Table 6.

In order to assess the statistical significance of our results, we repeated the experiments on 10 new random choices of subsets, using the hyperparameters selected on the original subset from Table 4 (except for learning rate, which is selected according to a different validation set for each subset). We then compared pairs of methods using a paired t-test, with p-values shown in Table 5. In particular, the results strengthen some of our findings, for instance, that  $\|\nabla f\|^2$  should be preferred to the gradient penalty on the loss when there is no data augmentation, and that combined upper+lower bound approaches tend to outperform the individual upper or lower bound strategies.

### A.2. Infinite MNIST

We provide more extensive results for the Infinite MNIST dataset in Table 7, in particular showing more regularization strategies, as well as results with or without data augmentation, marked with (\*). As in the case of CIFAR10, we use SGD with momentum (fixed to 0.9) for 500 epochs, with initial learning rates in [0.005; 0.05; 0.5], and divide the step-size by 2 every 40 epochs. The full hyper-parameter grid is given in Table 9.

As in the case of CIFAR10, we report statistical significance tests in Table 8 comparing pairs of methods based on 10 different random choices of subsets. In particular, the results confirm that weight decay with data augmentation alone tends to give weaker results than separate penalties, and that the combined penalty  $\|f\|_7^2 + \|f\|_8^2$ , which combines adversarial perturbations of two different types, outperforms each penalty taken by itself on a single type of perturbation, which emphasizes the benefit of considering perturbations of different natures, perhaps thanks to a tighter lower bound approximation of the RKHS norm. We note that  $\text{grad-}\ell_2(*)$  worked well on some subsets, but poorly on others due to training instabilities, possibly because of the selected hyperparameters which are quite large (and thus likely violate the approximation to the robust optimization objective).

### A.3. Protein homology detection

**Dataset description.** Our protein homology detection experiments consider the Structural Classification Of Proteins (SCOP) version 1.67 dataset (Murzin et al., 1995), filtered and split following the procedures of (Håndstad et al., 2007). Specifically, positive training samples are extracted from one superfamily from which one family is withheld to serve as positive test set, while negative sequences are chosen from outside of the target family’s hold and are randomly split into training and test samples in the same ratio as positive samples. This yields 102 superfamily classification tasks, which are generally very class-imbalanced. For each task, we sample 100 class-balanced training samples to use as training set. The positive samples are extended to 50 with Uniref50 using PSI-BLAST (Altschul et al., 1997) if they are fewer.

**Data augmentation procedure.** We consider in our experiments a discrete way of perturbing training samples to perform data augmentation. Specifically, for a given sequence, a perturbed sequence can be obtained by randomly changing some of the characters. Each character in the sequence is switched to a different one, randomly chosen from the alphabet, with some probability  $p$ . We fixed this probability to 0.1 throughout the experiments.

**Experimental details and significance tests.** In our experiments, we use the Adam optimization algorithm with a learning rate fixed to 0.01 (and  $\beta$  fixed to defaults (0.9, 0.999)), with a batch size of 100 for 300 epochs. The full hyper-parameter grid is given in Table 11. In addition to the average auROC50 scores reported in Table 3, we perform paired t-tests for

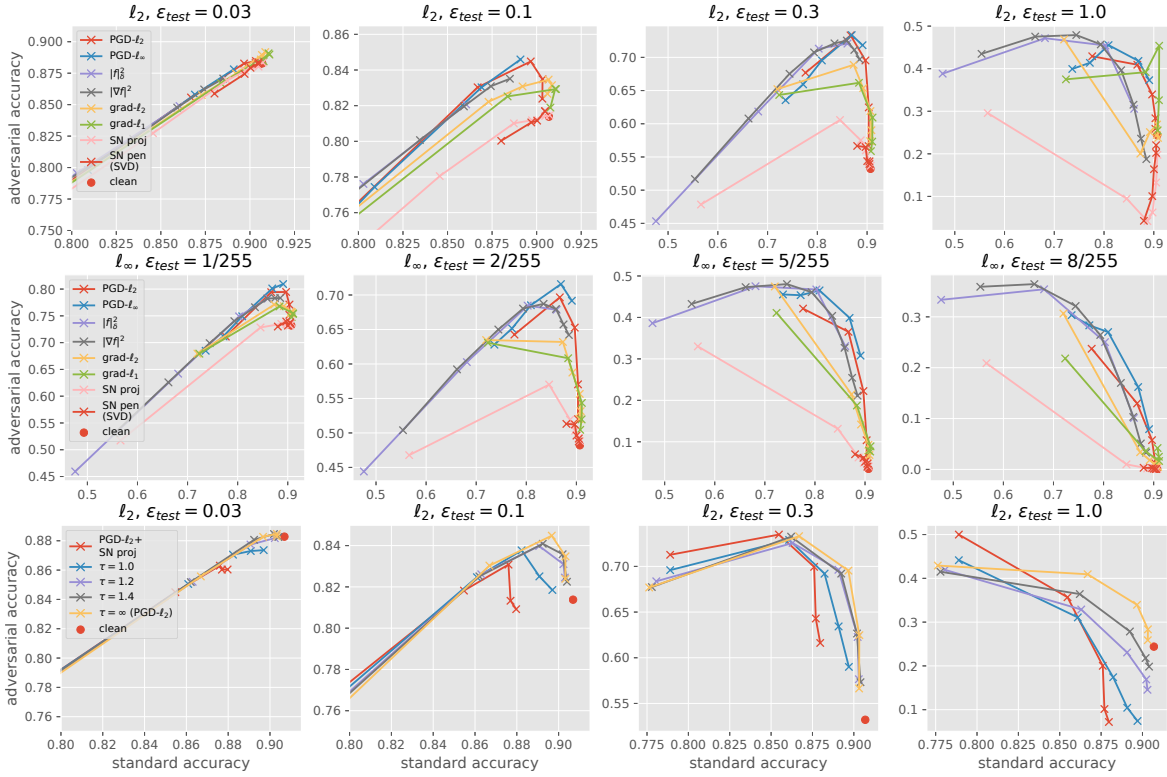


Figure 3. Robustness trade-off curves of different regularization methods for VGG11 on CIFAR10 (extended version of Figure 1). The plots show test accuracy vs adversarial test accuracy for  $l_2$ -bounded (top/bottom) or  $l_\infty$ -bounded (middle), 40-step PGD adversaries with a fixed  $\epsilon_{test}$ . Different points on a curve correspond to training with different regularization strengths. The regularization increases monotonically along a given curve, and the leftmost points correspond to the strongest regularization. The bottom plots consider PGD- $l_2$  + SN projection, with different fixed values of the constraint radius  $\tau$ , for varying  $\epsilon$  in PGD.

comparing pairs of methods in Table 10 in order to verify the significance of our findings. The results confirm that the adversarial perturbation penalty and its combination with spectral norm constraints tends to outperform the other approaches.

#### A.4. Robustness

Figure 3 extends Figure 1 from Section 4.2 to show more methods, adversary strengths, and different geometries. For combined (PGD- $l_2$  + SN projection) approaches, we can see that stronger constraints (*i.e.*, smaller  $\tau$ ) tend to reduce standard accuracy, likely because it prevents a good fit of the data, but can provide better robustness to strong adversaries ( $\epsilon_{test} = 1$ ). We can see that using the right metric in PGD indeed helps against an  $l_\infty$  adversary, nevertheless controlling global stability through the RKHS norm as in the  $\|f\|_\delta^2$  and  $\|\nabla f\|^2$  penalties can still provide some robustness against such adversaries, even with large  $\epsilon_{test}$ . For gradient penalties, we find that the different geometries behave quite similarly, which may suggest that more appropriate optimization algorithms than SGD could be needed to better accommodate the non-smooth case of  $l_1/l_\infty$ , or perhaps that both algorithms are actually controlling the same notion of complexity on this dataset.

Table 4. Regularization on CIFAR10 with 1 000 or 5 000 examples for VGG-11 and ResNet-18. Extended version of Table 1. Each entry shows the test accuracy with/without data augmentation when all hyper-parameters are optimized on a validation set of size 10 000 (a) or 1 000 (b), and for the epoch with highest validation accuracy, evaluating every 10 epochs (similar to early stopping).

(a) 10k examples in validation set

| Method                     | 1k VGG-11                   | 1k ResNet-18                | 5k VGG-11            | 5k ResNet-18         |
|----------------------------|-----------------------------|-----------------------------|----------------------|----------------------|
| No weight decay            | 50.70 / 43.75               | 45.23 / 37.12               | 72.49 / 58.35        | 72.72 / 54.12        |
| Weight decay               | 51.32 / 43.95               | 44.85 / 37.09               | 72.80 / 58.56        | 73.06 / 53.33        |
| SN penalty (PI)            | 54.64 / 45.06               | 47.01 / 39.63               | 74.03 / 62.45        | 74.79 / 54.04        |
| SN penalty (SVD)           | 53.44 / 46.06               | 47.26 / 37.94               | 74.53 / 62.93        | 75.59 / 54.98        |
| SN projection              | 54.14 / <b>46.70</b>        | 47.12 / 37.28               | 75.14 / 63.81        | 76.23 / 55.60        |
| VAT                        | 50.88 / 43.36               | 47.47 / 42.82               | 72.91 / 58.78        | 71.56 / 55.93        |
| PGD- $\ell_2$              | 51.25 / 44.40               | 45.80 / 41.87               | 73.18 / 58.98        | 72.53 / 55.92        |
| PGD- $\ell_\infty$         | 51.17 / 43.07               | 45.31 / 39.66               | 73.05 / 57.82        | 72.75 / 55.14        |
| grad- $\ell_2$             | <b>55.19</b> / 43.88        | <b>49.30</b> / <b>44.65</b> | <b>75.38</b> / 59.20 | 75.22 / 55.36        |
| grad- $\ell_1$             | 54.88 / 44.74               | <b>49.06</b> / 42.63        | <b>75.25</b> / 59.39 | 74.48 / 56.19        |
| $\ f\ _\delta^2$ penalty   | 51.41 / 45.07               | 48.73 / 43.72               | 72.98 / 61.45        | 72.78 / 56.50        |
| $\ \nabla f\ ^2$ penalty   | 54.80 / 46.37               | <b>48.99</b> / <b>44.97</b> | 73.90 / 60.17        | 73.83 / <b>57.92</b> |
| PGD- $\ell_2$ + SN proj    | 54.19 / <b>46.66</b>        | 47.47 / 41.25               | 74.61 / <b>64.50</b> | <b>77.19</b> / 57.43 |
| grad- $\ell_2$ + SN proj   | <b>55.32</b> / <b>46.88</b> | 48.73 / 42.78               | 75.11 / 63.54        | <b>77.73</b> / 57.09 |
| $\ f\ _\delta^2$ + SN proj | 54.02 / <b>46.72</b>        | 48.12 / 43.56               | 74.55 / <b>64.33</b> | 75.64 / <b>59.03</b> |
| $\ \nabla f\ ^2$ + SN proj | <b>55.24</b> / <b>46.80</b> | <b>49.06</b> / <b>44.92</b> | 72.31 / 63.74        | 72.24 / 57.56        |

(b) 1k examples in validation set

| Method                     | 1k VGG-11            | 1k ResNet-18         | 5k VGG-11            | 5k ResNet-18                |
|----------------------------|----------------------|----------------------|----------------------|-----------------------------|
| No weight decay            | 51.32 / 43.42        | 45.00 / 37.00        | 72.64 / 57.88        | 72.71 / 53.80               |
| Weight decay               | 51.04 / 43.42        | 44.66 / 36.77        | 72.68 / 57.59        | 72.25 / 54.16               |
| SN penalty (PI)            | 54.60 / 44.20        | 46.39 / 38.86        | 72.99 / 62.49        | 74.72 / 53.65               |
| SN penalty (SVD)           | 53.76 / 44.79        | 47.31 / 37.92        | 74.05 / 63.34        | 75.73 / 54.65               |
| SN projection              | 52.86 / <b>46.49</b> | 47.05 / 37.28        | 74.18 / <b>63.70</b> | 75.91 / 54.43               |
| VAT                        | 50.90 / 43.99        | 47.35 / 42.91        | 72.95 / 57.64        | 71.91 / 55.22               |
| PGD- $\ell_2$              | 50.95 / 43.26        | 45.77 / 41.71        | 72.71 / 57.68        | 72.87 / 54.17               |
| PGD- $\ell_\infty$         | 51.16 / 43.16        | 45.67 / 39.77        | 73.64 / 58.02        | 72.99 / 53.95               |
| grad- $\ell_2$             | <b>55.40</b> / 43.57 | 47.86 / <b>44.65</b> | <b>75.44</b> / 58.33 | 74.83 / 55.43               |
| grad- $\ell_1$             | 54.53 / 43.04        | <b>48.75</b> / 42.21 | <b>75.28</b> / 58.19 | 74.28 / 54.02               |
| $\ f\ _M^2$ penalty        | 51.00 / 44.67        | 48.57 / 44.30        | 72.76 / 60.55        | 72.75 / 56.49               |
| $\ \nabla f\ ^2$ penalty   | 54.68 / 46.10        | 48.53 / <b>45.21</b> | 73.83 / 60.36        | 73.30 / <b>57.46</b>        |
| PGD- $\ell_2$ + SN proj    | 53.85 / <b>46.79</b> | 46.48 / 40.95        | 74.79 / 63.37        | <b>76.28</b> / <b>57.43</b> |
| grad- $\ell_2$ + SN proj   | <b>55.28</b> / 45.11 | 48.42 / 41.93        | 75.17 / 63.45        | <b>77.24</b> / 56.18        |
| $\ f\ _M^2$ + SN proj      | 54.00 / 45.14        | 47.12 / 41.86        | 74.54 / <b>63.94</b> | 75.25 / <b>57.94</b>        |
| $\ \nabla f\ ^2$ + SN proj | <b>55.21</b> / 45.68 | <b>49.03</b> / 43.58 | 71.92 / 63.47        | 71.83 / 56.06               |

Table 5. Paired t-tests comparing pairs of methods, on 10 different random choices of subsets of CIFAR10. Each cell shows the p-value of the corresponding test, both with (left) and without (right) data augmentation. We only show p-values smaller than 0.05. Hyperparameters are fixed to the ones obtained for the results in Table 1 (selected on a different choice of subset), except for the learning rate which is tuned on a separate validation set for each choice of subset.

| Test  | 1k VGG-11 |       | 1k ResNet-18 |       | 5k VGG-11 |       | 5k ResNet-18 |       |
|---|-----------|-------|--------------|-------|-----------|-------|--------------|-------|
| SN projection $\succ$ Weight decay                  | 1e-04     | 1e-03 | -            | -     | 3e-06     | 1e-08 | 9e-07        | 4e-04 |
| grad- $\ell_2$ $\succ$ Weight decay                 | 4e-09     | -     | 2e-04        | 5e-05 | 7e-08     | 1e-04 | 5e-06        | -     |
| $\ \nabla f\ ^2$ $\succ$ Weight decay               | 1e-08     | 2e-07 | 1e-05        | 3e-07 | 3e-04     | 5e-07 | 7e-03        | 1e-06 |
| $\ \nabla f\ ^2$ $\succ$ grad- $\ell_2$             | -         | 3e-08 | 2e-02        | 2e-06 | -         | 6e-05 | -            | 4e-05 |
| grad- $\ell_2$ $\succ$ $\ \nabla f\ ^2$             | 2e-02     | -     | -            | -     | 2e-05     | -     | 7e-04        | -     |
| grad- $\ell_2$ + SN proj $\succ$ grad- $\ell_2$     | -         | 9e-03 | -            | -     | -         | 5e-07 | 9e-06        | 2e-04 |
| $\ \nabla f\ ^2$ + SN proj $\succ$ $\ \nabla f\ ^2$ | -         | -     | -            | 1e-02 | -         | 2e-06 | -            | -     |

Table 6. List of hyper-parameters used for each method on CIFAR10. For each method, we additionally consider a learning rate parameter in [0.003; 0.01; 0.03; 0.1]. For combined penalties, the sets of hyperparameters are listed in the same order as in the first column (*i.e.*, the choices of constraint radius are given last).

| Method                           | Parameter grid                                      |
|----------------------------------|---|
| No weight decay                  | -   |
| Weight decay                     | [0; 0.0001; 0.0002; 0.0004; 0.0008; 0.001; 0.002]   |
| SN penalty (PI)                  | [0.001; 0.003; 0.01; 0.03; 0.1; 0.3]                |
| SN penalty (SVD)                 | [0.001; 0.003; 0.01; 0.03; 0.1; 0.3]                |
| SN projection                    | [0.5; 0.6; 0.8; 1.0; 1.2; 1.4]                      |
| $\ f\ _\delta^2$ penalty         | [0.001; 0.003; 0.01; 0.03; 0.1]                     |
| $\ \nabla f\ ^2$ penalty         | [0.00003; 0.0001; 0.0003; 0.001; 0.003; 0.01; 0.03] |
| VAT                              | [0.1; 0.3; 1.0; 3.0]                                |
| PGD- $\ell_2$                    | [0.003; 0.01; 0.03; 0.1; 0.3; 1.0]                  |
| PGD- $\ell_\infty$               | [0.001; 0.003; 0.01; 0.03; 0.1; 0.3]                |
| grad- $\ell_1$                   | [0.0001; 0.0003; 0.001; 0.003; 0.01; 0.03]          |
| grad- $\ell_2$                   | [0.001; 0.003; 0.01; 0.03; 0.1; 0.3; 1.0; 3.0]      |
| PGD- $\ell_2$ + SN projection    | [0.003; 0.01; 0.03; 0.1] $\times$ [0.6; 1.0; 1.4]   |
| grad- $\ell_2$ + SN projection   | [0.003; 0.01; 0.03; 0.1] $\times$ [0.6; 1.0; 1.4]   |
| $\ f\ _\delta^2$ + SN projection | [0.003; 0.01; 0.03] $\times$ [0.6; 1.0; 1.4]        |
| $\ \nabla f\ ^2$ + SN projection | [0.001; 0.01; 0.1] $\times$ [0.6; 1.0; 1.4]         |



Table 7. Test accuracies on subsets of MNIST using deformations from Infinite MNIST. Extended version of Table 2. (\*) indicates that random deformations were included as training examples (*i.e.*, data augmentation), while  $\|f\|_\tau^2$  and  $\|D_\tau f\|^2$  use them as part of the regularization penalty. As in Table 4, we show results obtained using a validation set of size 10 000 (a) and 1 000 (b).

| (a) 10k examples in validation set            |              |              | (b) 1k examples in validation set             |              |              |
|---|--------------|--------------|---|--------------|--------------|
| Method  | 300 VGG      | 1k VGG       | Method  | 300 VGG      | 1k VGG       |
| Weight decay                                  | 89.32        | 94.08        | Weight decay                                  | 89.32        | 93.34        |
| Weight decay (*)                              | 92.41        | 95.64        | Weight decay (*)                              | 91.91        | 95.73        |
| SN projection                                 | 90.69        | 95.01        | SN projection                                 | 90.60        | 94.83        |
| SN projection (*)                             | 92.17        | 95.88        | SN projection (*)                             | 92.01        | 95.91        |
| grad- $\ell_2$                                | 93.63        | 96.67        | grad- $\ell_2$                                | 92.92        | 96.42        |
| grad- $\ell_2$ (*)                            | 95.05        | 97.48        | grad- $\ell_2$ (*)                            | <b>94.69</b> | <b>97.48</b> |
| $\ f\ _\delta^2$ penalty                      | 94.17        | 96.99        | $\ f\ _M^2$ penalty                           | 93.44        | 96.98        |
| $\ f\ _\delta^2$ penalty (*)                  | 94.86        | 97.40        | $\ f\ _M^2$ penalty (*)                       | 94.57        | 97.14        |
| $\ \nabla f\ ^2$ penalty                      | 94.08        | 96.82        | $\ \nabla f\ ^2$ penalty                      | 94.08        | 96.77        |
| $\ \nabla f\ ^2$ penalty (*)                  | 94.80        | 97.29        | $\ \nabla f\ ^2$ penalty (*)                  | 94.50        | 97.15        |
| $\ D_\tau f\ ^2$ penalty                      | 94.18        | 96.98        | $\ D_\tau f\ ^2$ penalty                      | 94.03        | 97.16        |
| $\ D_\tau f\ ^2$ penalty (*)                  | 94.91        | 97.29        | $\ D_\tau f\ ^2$ penalty (*)                  | 94.15        | 96.64        |
| $\ f\ _\tau^2$ penalty                        | 94.42        | 97.13        | $\ f\ _\tau^2$ penalty                        | 93.53        | 97.13        |
| $\ f\ _\tau^2$ penalty (*)                    | 94.83        | 97.25        | $\ f\ _\tau^2$ penalty (*)                    | <b>94.79</b> | 97.26        |
| $\ f\ _\tau^2 + \ \nabla f\ ^2$               | 94.75        | 97.40        | $\ f\ _\tau^2 + \ \nabla f\ ^2$               | <b>94.75</b> | 97.21        |
| $\ f\ _\tau^2 + \ \nabla f\ ^2$ (*)           | 95.14        | 97.44        | $\ f\ _\tau^2 + \ \nabla f\ ^2$ (*)           | 94.43        | <b>97.42</b> |
| $\ f\ _\tau^2 + \ f\ _\delta^2$               | 95.23        | <b>97.66</b> | $\ f\ _\tau^2 + \ f\ _M^2$                    | <b>95.15</b> | 97.27        |
| $\ f\ _\tau^2 + \ f\ _\delta^2$ (*)           | <b>95.53</b> | <b>97.56</b> | $\ f\ _\tau^2 + \ f\ _M^2$ (*)                | <b>95.20</b> | <b>97.49</b> |
| grad- $\ell_2$ + SN proj                      | 93.89        | 96.85        | grad- $\ell_2$ + SN proj                      | 93.44        | 96.81        |
| grad- $\ell_2$ + SN proj (*)                  | 95.15        | <b>97.80</b> | grad- $\ell_2$ + SN proj (*)                  | 94.05        | <b>97.60</b> |
| $\ f\ _\delta^2 +$ SN proj                    | 93.97        | 96.89        | $\ f\ _M^2 +$ SN proj                         | 93.97        | 96.61        |
| $\ f\ _\delta^2 +$ SN proj (*)                | 94.78        | 97.38        | $\ f\ _M^2 +$ SN proj (*)                     | <b>94.69</b> | 97.33        |
| $\ f\ _\tau^2 + \ \nabla f\ ^2 +$ SN proj     | 95.09        | 97.42        | $\ f\ _\tau^2 + \ \nabla f\ ^2 +$ SN proj     | <b>94.75</b> | 97.16        |
| $\ f\ _\tau^2 + \ \nabla f\ ^2 +$ SN proj (*) | 95.03        | 97.27        | $\ f\ _\tau^2 + \ \nabla f\ ^2 +$ SN proj (*) | <b>94.74</b> | 97.22        |
| $\ f\ _\tau^2 + \ f\ _\delta^2 +$ SN proj     | 95.20        | <b>97.60</b> | $\ f\ _\tau^2 + \ f\ _M^2 +$ SN proj          | <b>94.78</b> | <b>97.49</b> |
| $\ f\ _\tau^2 + \ f\ _\delta^2 +$ SN proj (*) | <b>95.40</b> | <b>97.77</b> | $\ f\ _\tau^2 + \ f\ _M^2 +$ SN proj (*)      | <b>95.17</b> | <b>97.64</b> |

Table 8. Paired t-tests comparing pairs of methods, on 10 different random choices of subsets of MNIST. Each cell shows the p-value of the corresponding test. We only show p-values smaller than 0.05. Hyperparameters are fixed to the ones obtained for the results in Table 2 (selected on a different choice of subset), except for the learning rate which is tuned on a separate validation set for each choice of subset.

| Test  | 300 VGG | 1k VGG |
|---|---------|--------|
| grad- $\ell_2$ (*) $\succ$ Weight decay (*)   | -       | 3e-11  |
| $\ f\ _\tau^2$ penalty $\succ$ Weight decay (*)                                       | 2e-08   | 2e-10  |
| $\ f\ _\tau^2 + \ f\ _\delta^2$ $\succ$ Weight decay (*)                              | 1e-08   | 2e-10  |
| $\ f\ _\tau^2 + \ f\ _\delta^2 +$ SN proj (*) $\succ$ grad- $\ell_2$ (*)              | -       | 1e-02  |
| grad- $\ell_2$ (*) $\succ$ $\ f\ _\tau^2 + \ f\ _\delta^2 +$ SN proj (*)              | -       | -      |
| $\ f\ _\tau^2 + \ f\ _\delta^2$ $\succ$ $\ f\ _\delta^2$ penalty                      | 1e-07   | 6e-09  |
| $\ f\ _\tau^2 + \ f\ _\delta^2$ $\succ$ $\ f\ _\tau^2$ penalty                        | 2e-06   | 6e-07  |
| $\ f\ _\tau^2 + \ f\ _\delta^2$ (*) $\succ$ $\ f\ _\tau^2 + \ f\ _\delta^2$           | 2e-03   | -      |
| $\ f\ _\tau^2 + \ f\ _\delta^2 +$ SN proj (*) $\succ$ $\ f\ _\tau^2 + \ f\ _\delta^2$ | 2e-03   | 2e-04  |

Table 9. List of hyper-parameters used for each method on Infinite MNIST. For each method, we additionally consider a learning rate parameter in  $[0.005; 0.05; 0.5]$ . For combined penalties, the sets of hyperparameters are listed in the same order as in the first column (e.g., the choices of constraint radius are given last).

| Method                                    | Grid   |
|---|--|
| Weight decay                              | [0; 0.00001; 0.00003; 0.0001; 0.0003; 0.001; 0.003; 0.01; 0.03; 0.1] |
| SN projection                             | [1.0; 1.2; 1.4; 1.6; 1.8]  |
| grad- $\ell_2$                            | [0.1; 0.3; 1.0; 3.0; 10.0]   |
| $\ f\ _\delta^2$ penalty                  | [0.1; 0.3; 1.0; 3.0]   |
| $\ \nabla f\ ^2$ penalty                  | [0.0003; 0.001; 0.003; 0.01; 0.03; 0.1; 0.3]                         |
| $\ D_\tau f\ ^2$ penalty                  | [0.003; 0.01; 0.03; 0.1; 0.3]  |
| $\ f\ _\tau^2$ penalty                    | [0.03; 0.1; 0.3; 1.0; 3.0]   |
| $\ f\ _\tau^2 + \ \nabla f\ ^2$           | [0.03; 0.1; 0.3; 1.0] $\times$ [0.003; 0.01; 0.03; 0.1]              |
| $\ f\ _\tau^2 + \ f\ _\delta^2$           | [0.1; 0.3; 1.0] $\times$ [0.03; 0.1]                                 |
| grad- $\ell_2$ + SN proj                  | [0.3; 1.0; 3.0; 10.0; 30.0] $\times$ [1.2; 1.6; 2.0]                 |
| $\ f\ _\delta^2$ + SN proj                | [0.03; 0.1] $\times$ [1.2; 1.6; 2.0]                                 |
| $\ f\ _\tau^2 + \ \nabla f\ ^2$ + SN proj | [0.03; 0.1; 0.3] $\times$ [0.01; 0.03; 0.1] $\times$ [1.2; 1.6; 2.0] |
| $\ f\ _\tau^2 + \ f\ _\delta^2$ + SN proj | [0.1; 0.3; 1.0] $\times$ [0.03; 0.1] $\times$ [1.2; 1.6; 2.0]        |

Table 10. Paired t-tests comparing pairs of methods on the 51 test datasets from the set of protein homology detection tasks. Each cell shows the p-value of the corresponding test. We only show p-values smaller than 0.05. We use the same hyperparameters as the ones obtained in the results of Table 3.

| Test  | No DA | DA    |
|---|-------|-------|
| SN proj $\succ$ Weight decay                                  | 1e-05 | 4e-05 |
| grad- $\ell_2$ $\succ$ Weight decay                           | 5e-05 | 5e-02 |
| $\ f\ _\delta^2$ $\succ$ Weight decay                         | 5e-06 | 3e-05 |
| $\ \nabla f\ ^2$ $\succ$ Weight decay                         | 9e-06 | 3e-03 |
| $\ f\ _\delta^2$ $\succ$ grad- $\ell_2$                       | -     | 4e-03 |
| $\ \nabla f\ ^2$ $\succ$ grad- $\ell_2$                       | -     | -     |
| grad- $\ell_2$ + SN proj $\succ$ grad- $\ell_2$               | -     | 1e-03 |
| $\ f\ _\delta^2$ + SN proj $\succ$ $\ f\ _\delta^2$           | 3e-03 | 5e-02 |
| $\ \nabla f\ ^2$ + SN proj $\succ$ $\ \nabla f\ ^2$           | -     | -     |
| $\ f\ _\delta^2$ + SN proj $\succ$ $\ \nabla f\ ^2$ + SN proj | 8e-05 | -     |

Table 11. List of hyper-parameters used for each method on protein homology detection datasets. For combined penalties, the hyperparameters are the cross-products of each individual method.

| Method           | Parameter grid                        |
|------------------|---------------------------------------|
| No weight decay  | -                                     |
| Weight decay     | [0; 0.01; 0.001; 0.0001; 0.00001]     |
| SN proj          | [10; 1.0; 0.1]                        |
| PGD- $\ell_2$    | [100.0; 10.0; 1.0; 0.1]               |
| grad- $\ell_2$   | [100.0; 10.0; 1.0; 0.1; 0.01, 0.001]  |
| $\ f\ _\delta^2$ | [10.0; 1.0; 0.1]                      |
| $\ \nabla f\ ^2$ | [10.0; 1.0; 0.1; 0.01; 0.001; 0.0001] |

## B. Details on Deformation Stability Penalties

This section provides more details on the deformation stability penalties mentioned in Section 2.2, and the practical versions we use in our experiments on the Infinite MNIST dataset (Loosli et al., 2007).

**Stability to deformations.** We begin by providing some background on deformation stability, recalling that these can provide new lower bound penalties as explained in Section 2.2. Viewing an element  $x \in \mathcal{X}$  as a signal  $x(u)$ , where  $u$  denotes the location (e.g. a two-dimensional vector for images), we denote by  $x_\tau$  a deformed version of  $x$  given by  $x_\tau(u) = x(u - \tau(u))$ , where  $\tau$  is a diffeomorphism. The deformation stability bounds of Bietti & Mairal (2019) take the form:

$$\|\Phi(x_\tau) - \Phi(x)\|_{\mathcal{H}} \leq (C_1\|\tau\|_\infty + C_2\|\nabla\tau\|_\infty)\|x\|, \quad (14)$$

where  $\nabla\tau(u)$  is the Jacobian of  $\tau$  at location  $u$ . Here,  $C_1$  controls translation invariance and typically decreases with the total amount of pooling (i.e., translation invariance more or less corresponds to the resolution at the final layer), while  $C_2$  controls stability to deformations (note that  $\nabla\tau = 0$  for translations) and is typically smaller when using small patches. We note that the bounds assume linear pooling layers with a certain spatial decay, adapted to the resolution of the current layer; our experiments on Infinite MNIST with deformation stability penalties thus use average pooling layers on 2x2 neighborhoods.

**Adversarial deformation penalty.** We can obtain lower bound penalties by exploiting the above stability bounds in a similar manner to the adversarial perturbation penalty introduced in Section 2.2. In particular, assuming a scalar-valued convolutional network  $f$ :

$$\|f\|_\tau^2 := \sup_{x \in \mathcal{X}, \tau \in \mathcal{T}} (f(x_\tau) - f(x))^2 \quad (15)$$

where  $\mathcal{T}$  is a collection of diffeomorphisms. When the diffeomorphisms in  $\mathcal{T}$  have bounded norm  $\|\tau\|_\infty$  and Jacobian norm  $\|\nabla\tau\|_\infty$ , and assuming  $\mathcal{X}$  (or, in practice, the training data) is bounded, the stability bound 14 ensures that the set  $U_\mathcal{T} = \{\Phi(x_\tau) - \Phi(x) : x \in \mathcal{X}, \tau \in \mathcal{T}\}$  is included in an RKHS ball with some radius  $r$ , so that  $\|f\|_\tau$  is a lower bound on  $r\|f\|_{\mathcal{H}}$ .

**Tangent gradient penalty.** We also consider the following gradient penalty along tangent vectors, which provides an approximation of the above adversarial penalty when considering small, parameterized deformations, and recovers the tangent propagation strategy of Simard et al. (1998):

$$\|D_\tau f\|^2 := \sup_{x \in \mathcal{X}} \|\partial_\alpha f(x + \sum_i \alpha_i t_{x,i})\|^2, \quad (16)$$

where  $\{t_{x,i}\}_{i=1,\dots,q}$  are tangent vectors at  $x$  obtained from a given set of deformations. To see the link with the adversarial deformation penalty 15, consider for simplicity a single deformation,  $\mathcal{T} = \{\tau_0\}$ . For small  $\alpha$ , we have

$$x_{\alpha\tau_0} \approx x + \alpha t_x, \quad \text{where } t_x(u) = \tau_0(u) \cdot \nabla x(u),$$

where  $t_x$  denotes the tangent vector of the deformation manifold  $\{\alpha\tau_0 : \alpha\}$  at  $\alpha = 0$  (Simard et al., 1998). Then,

$$f(x_{\alpha\tau_0}) - f(x) \approx \alpha \partial_\alpha f(x + \alpha t_x) = \alpha \langle \nabla f(x), t_x \rangle.$$

In this case, denoting  $\alpha\mathcal{T} = \{\alpha\tau_0\}$ , we have

$$\sup_{x \in \mathcal{X}, \tau \in \alpha\mathcal{T}} (f(x_\tau) - f(x))^2 \approx \alpha^2 \sup_{x \in \mathcal{X}} |\partial_\alpha f(x + \alpha t_x)|^2,$$

so that when  $\alpha$  is small, the adversarial penalty can be approximated by  $\alpha\|D_\tau f\|$  (note that using  $\alpha\mathcal{T}$  instead of  $\mathcal{T}$  in the adversarial penalty would also yield a scaling by  $\alpha$ , since the stability bounds imply  $\alpha$  times smaller perturbations in the RKHS).

**Practical implementations on Infinite MNIST.** In our experiments on Infinite MNIST, we compute  $\|f\|_\tau^2$  by considering 32 random transformations of each digit in a mini-batch of training examples, and taking the maximum over both the example and the transformation. We do this separately for each class, as for the other lower bound penalties  $\|f\|_\delta^2$  and  $\|\nabla f\|^2$ . For  $\|D_\tau f\|^2$ , we take  $\{t_{x,i}\}_{i=1,\dots,q}$  with  $q = 30$  to be tangent vectors given by random diffeomorphisms from Infinite MNIST around each example  $x$ .

### C. Details on Optimization with Spectral Norms

This section details our optimization approach presented in Section 2.3 for learning with spectral norm constraints. In particular, we rely on a *continuation* approach, decreasing the size of the ball constraints during training, towards a final value  $\tau$ . The method is presented in Algorithm 1. We use an exponentially decreasing schedule for  $\tau$ , and take  $\kappa$  to be 2 epochs for regularization, and 50 epochs for robustness. In the context of convolutional networks, we simply consider the SVD of a reshaped filter matrix, but we note that alternative approaches based on the singular values of the full convolutional operation may also be used (Sedghi et al., 2019).

---

**Algorithm 1** Stochastic projected gradient with continuation

---

```

Input:  $\tau, \kappa$ , step-sizes  $\eta_t$ 
for  $t = 1, \dots$  do
  Sample mini-batch and compute gradients of the loss w.r.t. each  $W^l$ , denoted  $G_t^l$ 
   $\tau_t = \tau(1 + \exp(-\frac{t}{\kappa}))$ 
  for  $l = 1, \dots, L$  do
     $\tilde{W}_t^l := W_t^l - \eta_t G_t^l$ 
    Compute SVD:  $\tilde{W}_t^l = U \text{diag}(\sigma) V^T$ 
     $\hat{\sigma} := \text{proj}_{\|\cdot\|_\infty \leq \tau_t}(\sigma)$ 
     $W_{t+1}^l := U \text{diag}(\hat{\sigma}) V^T$ 
  end for
end for

```

---

### D. Extensions to Non-Euclidian Geometries

The kernel approach from previous sections is well-suited for input spaces  $\mathcal{X}$  equipped with the Euclidian distance, thanks to the non-expansiveness property (2) of the kernel mapping. In the case of linear models, this kernel approach corresponds to using  $\ell_2$ -regularization by taking a linear kernel. However, other forms of regularization and geometries can often be useful, for example to encourage sparsity with an  $\ell_1$  regularizer. Such a regularization approach presents tight links with robustness to  $\ell_\infty$  perturbations on input data, thanks to the duality relation  $\|w\|_1 = \sup_{\|u\|_\infty} \langle w, u \rangle$  (see Xu et al., 2009a).

In the context of deep networks, we can leverage such insights to obtain new regularizers, expressed in the same variational form as the lower bounds in Section 2.2, but with different geometries on  $\mathcal{X}$ . For  $\ell_\infty$  perturbations, we obtain

$$\sup_{x, y \in \mathcal{X}} \frac{f(x) - f(y)}{\|x - y\|_\infty} \geq \sup_{x \in \mathcal{X}} \|\nabla f(x)\|_1. \quad (17)$$

The Lipschitz regularizer (l.h.s.) can also be taken in an adversarial perturbation form, with  $\ell_\infty$ -bounded perturbations  $\|\delta\|_\infty \leq \epsilon$ . When considering the corresponding robust optimization problem

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^n \sup_{\|\delta\|_\infty \leq \epsilon} \ell(y_i, f_{\theta}(x_i + \delta)), \quad (18)$$

we may consider the PGD approach of Madry et al. (2018), or the associated gradient penalty approach with the  $\ell_1$  norm, which is a good approximation when  $\epsilon$  is small (Lyu et al., 2015; Simon-Gabriel et al., 2019).

As most visible in the gradient  $\ell_1$ -norm in (17), these penalties encourage some sparsity in the gradients of  $f$ , which is a reasonable prior for regularization on images, for instance, where we might only want predictions to change based on few salient pixel regions. This can lead to gains in interpretability, as observed by Tsipras et al. (2019).

We note that in the case of linear models, our robust margin bound of Section 3.1 can be adapted to  $\ell_\infty$ -perturbations, by leveraging Rademacher complexity bounds for  $\ell_1$ -constrained models (Kakade et al., 2009). Obtaining similar bounds for neural networks would be interesting but goes beyond the scope of this paper.

### E. Details on Generalization Guarantees

This section presents the proof of Proposition 1, which relies on standard tools from statistical learning theory (e.g., Boucheron et al., 2005).



### E.1. Proof of Proposition 1

*Proof.* Assume for now that  $\gamma$  is fixed in advance, and let  $\mathcal{F}_\lambda := \{f \in \mathcal{H} : \|f\|_{\mathcal{H}} \leq \lambda\}$ . Note that for all  $f \in \mathcal{F}_\lambda$  we have

$$\text{err}_{\mathcal{D}}(f, \epsilon) = P(\exists \|\delta\| \leq \epsilon : yf(x + \delta) < 0) \leq P(yf(x) < \lambda\epsilon) =: L^{\lambda\epsilon}(f),$$

since  $\|f\|_{\mathcal{H}} \leq \lambda$  is an upper bound on the Lipschitz constant of  $f$ . Consider the function

$$\phi(x) = \begin{cases} 0, & \text{if } x \leq -\gamma - \lambda\epsilon \\ 1, & \text{if } x \geq -\lambda\epsilon \\ 1 + (x + \lambda\epsilon)/\gamma, & \text{otherwise.} \end{cases}$$

Defining  $A(f) = \mathbb{E} \phi(-yf(x)) \geq L^{\lambda\epsilon}(f)$  and  $A_n(f) = \frac{1}{n} \sum_{i=1}^n \phi(-y_i f(x_i)) \leq L_n^{\lambda\epsilon+\gamma}(f)$ , and noting that  $\phi$  is upper bounded by 1 and  $1/\gamma$  Lipschitz, we can apply similar arguments to (Boucheron et al., 2005, Theorem 4.1) to obtain, with probability  $1 - \delta$ ,

$$L^{\lambda\epsilon}(f) \leq L_n^{\lambda\epsilon+\gamma}(f) + O\left(\frac{1}{\gamma} R_n(\mathcal{F}_\lambda) + \sqrt{\frac{\log 1/\delta}{n}}\right),$$

where  $R_n(\mathcal{F}_\lambda)$  denotes the empirical Rademacher complexity of  $\mathcal{F}_\lambda$  on the dataset  $\{(x_i, y_i)\}_{i=1, \dots, n}$ . Standard upper bounds on empirical Rademacher complexity of kernel classes with bounded RKHS norm yield the following bound

$$\text{err}_{\mathcal{D}}(f, \epsilon) \leq L_n^{\lambda\epsilon+\gamma}(f) + O\left(\frac{\lambda}{\gamma\sqrt{n}} \sqrt{\frac{1}{n} \sum_{i=1}^n K(x_i, x_i)} + \sqrt{\frac{\log 1/\delta}{n}}\right).$$

Note that the bound is still valid with  $\gamma' \geq \gamma$  instead of  $\gamma$  in the first term of the r.h.s., since  $L_n^\gamma(f)$  is non-decreasing as a function of  $\gamma$ .

In order to establish the final bound, we instantiate the previous bound for values  $\lambda_i = 2^i$  and  $\gamma_j = 2^{-j}$ . Defining  $\delta_{i,j} = \frac{\delta}{(1+4i^2) \cdot (1+4j^2)}$ , we have that w.p.  $1 - \delta_{i,j}$ , for all  $f \in \mathcal{F}_{\lambda_i}$  and all  $\gamma \geq \gamma_j$ ,

$$\text{err}_{\mathcal{D}}(f, \epsilon) \leq L_n^{\lambda_i\epsilon+\gamma}(f) + O\left(\frac{\lambda_i}{\gamma_j\sqrt{n}} \sqrt{\frac{1}{n} \sum_{i=1}^n K(x_i, x_i)} + \sqrt{\frac{\log 1/\delta_{i,j}}{n}}\right). \quad (19)$$

By a union bound, this event holds jointly for all integers  $i, j$  w.p. greater than  $1 - \delta$ , since  $\sum_{i,j} \delta_{i,j} \leq \delta$ . Now consider an arbitrary  $f \in \mathcal{H}$  and  $\gamma > 0$  and let  $i = \lceil \log_2 \|f\|_{\mathcal{H}} \rceil$  and  $j = \lceil \log_2(1/\gamma) \rceil$ . We have

$$\begin{aligned} \lambda_i &\leq 2\|f\|_{\mathcal{H}} \\ \frac{1}{\gamma_j} &\leq \frac{2}{\gamma} \\ \log(1/\delta_{i,j}) &\leq \log(C(\|f\|_{\mathcal{H}}, \gamma)/\delta), \end{aligned}$$

with  $C(\|f\|_{\mathcal{H}}, \gamma) := (1 + 4(\log_2 \|f\|_{\mathcal{H}})^2) \cdot (1 + 4(\log_2(1/\gamma))^2)$ . Applying this to the bound in (19) yields the desired result.  $\square$