
Adversarial Learning Guarantees for Linear Hypotheses and Neural Networks

Pranjal Awasthi¹ Natalie S. Frank² Mehryar Mohri³

Abstract

Adversarial or test time robustness measures the susceptibility of a classifier to perturbations to the test input. While there has been a flurry of recent work on designing defenses against such perturbations, the theory of adversarial robustness is not well understood. In order to make progress on this, we focus on the problem of understanding generalization in adversarial settings, via the lens of Rademacher complexity. We give upper and lower bounds for the adversarial empirical Rademacher complexity of linear hypotheses with adversarial perturbations measured in l_r -norm for an arbitrary $r \geq 1$. We then extend our analysis to provide Rademacher complexity lower and upper bounds for a single ReLU unit. Finally, we give adversarial Rademacher complexity bounds for feed-forward neural networks with one hidden layer.

1. Introduction

Robustness is a key requirement when designing machine learning models and comes in various forms such as robustness to training set corruptions, missing feature values, and model misspecification. In recent years, requiring robustness to *adversarial* or *test time* perturbations has become a key requirement. Starting with the work of [Szegedy et al. \(2014\)](#) it has now been well established that deep neural networks trained via standard gradient descent based algorithms are highly susceptible to imperceptible corruptions to the input at test time ([Goodfellow et al., 2014](#); [Chen et al., 2017](#); [Eykholt et al., 2018](#); [Carlini & Wagner, 2018](#)). This has led to a proliferation of work aimed at designing classifiers robust to such perturbations ([Madry et al., 2017](#);

[Gowal et al., 2018](#); [2019](#); [Schott et al., 2018](#)) and works aimed at designing more sophisticated attacks to break such classifiers ([Athalye et al., 2018](#); [Carlini & Wagner, 2017](#); [Sharma & Chen, 2017](#))

While the above works have made significant progress in designing practical defenses, theoretical aspects of adversarial robustness are currently poorly understood unlike other notions of training set corruptions that have been widely studied in both the statistics and the computer science communities ([Huber, 2011](#); [Kearns & Li, 1993](#); [Kearns et al., 1994](#)). Theoretical understanding of adversarial robustness presents three main challenges. The first is a computational one since even checking the robustness of a given model at a given test input is an NP-hard problem ([Awasthi et al., 2019](#)). This has been explored in recent works that construct specific instances of learning problems where standard non-robust learning can be done efficiently, but learning a robust classifier becomes computationally hard ([Bubeck et al., 2018b;a](#); [Nakkiran, 2019](#); [Degwekar et al., 2019](#)). The second challenge concerns whether achieving adversarial robustness requires one to compromise on standard accuracy. Recent works have shown specific instance where this tradeoff is inherent ([Tsipras et al., 2018](#); [Raghunathan et al., 2019](#)).

Finally, the third challenge, the main focus of this work, is the question of what quantity governs generalization in adversarial settings, and how generalization in adversarial settings compares to its non-adversarial counterpart. The recent work of [Schmidt et al. \(2018\)](#) has shown, via specific constructions, that in some scenarios achieving adversarial generalization requires more data as compared to adversarial generalization. Furthermore, the work of [Montasser et al. \(2019\)](#) casts a shadow of doubt on the use of classical quantities such as the VC-dimension of explain generalization in adversarial settings.

However, generalization for function classes of infinite VC-dimension (such as SVMs with a Gaussian kernel) can be explained via margin-based bounds. Characterizing the Rademacher complexity of the function class is essential in these estimates. In a similar vein, we believe that providing non-trivial bounds on the adversarial Rademacher complexity can help shed light on when generalization is possible in

¹Google Research and Rutgers University ²Courant Institute of Math. Sciences ³Google Research and Courant Institute of Math. Sciences. Correspondence to: Natalie S. Frank <nf1066@nyu.edu>.

adversarial settings via similar margin-based bounds. The difficulty is that current bounds on adversarial Rademacher complexity are too loose and vacuous in many settings. This is the barrier that we aim to overcome in this work.

In order to make progress on the mystery of adversarial generalization, a recent line of work (Khim & Loh, 2018; Yin et al., 2019) aims to study the notion of Rademacher complexity for various function classes in the adversarial settings. Focusing mainly on the case of linear models, these works aim to quantify the additional overhead in sample complexity that is incurred when requiring adversarial generalization. Extending the ideas to the case of more general neural networks becomes more challenging and as a result there works instead bound the Rademacher complexity in terms of the Rademacher complexity of an appropriate surrogate. In this work we extend this line of work along several directions.

Our Contributions. We provide a general analysis of the adversarial Rademacher complexity of linear models that holds for perturbations measured in any ℓ_p norm. This extends the prior work of Yin et al. (2019) that applies only to ℓ_∞ adversarial perturbation and provided a finer analysis of linear models as compared to the work of Khim & Loh (2018).

As a consequence of our analysis, we provide a sharp characterization of when the adversarial Rademacher complexity suffers from an additional dimension dependent term as compared to its non-adversarial counterpart. This has algorithmic implications for designing appropriate regularizers for adversarial learning of linear models. As an additional byproduct, we are able to provide improved Rademacher complexity bounds for linear classifiers, even in non-adversarial scenarios!

As a next step towards understanding neural networks, we then extend our analysis to provide data dependent upper and lower bounds on the adversarial Rademacher complexity of a single ReLU unit.

Finally, we provide upper bounds on the adversarial Rademacher complexity of one hidden layer neural networks. In contrast with prior works (Yin et al., 2019; Khim & Loh, 2018), our bounds directly apply to the original network as opposed to a surrogate. Our bounds for neural networks come in two forms. We first provide a general upper bound that applies to any neural network with Lipschitz activations. This bound as a dependence on the underlying dimensionality of the input data. Next, we provide a finer data dependent upper bound that is related to the ϵ -adversarial growth function of the data, a quantity we introduce in this work.

Comparison with Prior Work Yin et al. (2019) and Khim and Loh (2018) previously studied the adversarial

Rademacher complexity of linear classifiers and neural networks. Our work adds to this line of research in multiple ways.

Yin et al. (2019) analyzed the adversarial Rademacher complexity of linear models when perturbations are measured in ℓ_∞ norm. They show that in this case the adversarial Rademacher complexity of the loss class is bounded above and below by the sum of its non-adversarial counterpart and a dimension dependent term. Our result is a strict generalization of (Yin et al., 2019) because we provide the analysis of adversarial Rademacher complexity when the perturbations are measured in any general ℓ_r norm.

The recent work of Khim and Loh (2018) also studies the adversarial Rademacher complexity of classes of linear models constrained by the 2-norm under general ℓ_r perturbations. They noted that the upper bound on the adversarial Rademacher complexity involves a constant that depends on r . While the bounds are qualitatively similar, they did not derive the relationship with the dimension d nor did they give a lower bound involving this term. ((Yin et al., 2019) noted the relationship to d only for $r = \infty$.) Further, we consider models bounded by any ℓ_p norm in addition to the 2-norm, and this consideration turns out to have important implications for generalization and model selection. In the process, we improve upon the existing classical analysis of (non-adversarial) Rademacher complexity of linear models, which is of independent interest.

For the case of neural networks, both the works of Yin et al. (2019) and Khim and Loh (2018) replace the adversarial loss defined as $\min_{\mathbf{x}': \|\mathbf{x}' - \mathbf{x}\| \leq \epsilon} \phi(yf(\mathbf{x}'))$, by a surrogate upper bound and analyze the resulting Rademacher complexity of the surrogate. In general, these bounds on the surrogate might not lead to meaningful generalization bounds on the original adversarial loss.

In the work of Yin et al. (2019) the surrogate is chosen to be an upper bound on the adversarial loss based on a semi-definite programming (SDP) based relaxation. However, the relaxation is quite weak and is obtained via taking the worst possible \mathbf{x} in the domain of the input. This approach greatly overestimates the adversarial loss.

In the work of (Khim & Loh, 2018) the surrogate is based on the adversarial loss of another neural network that is derived from the original one via a tree based decomposition. A key issue in adversarial robustness is the analysis of the optimization problem: $\sup_{\mathbf{s}} f(\mathbf{x} + \mathbf{s})$, where \mathbf{x} is a given input and \mathbf{s} is a perturbation chosen in a small ball. Understanding the structure of the optimal \mathbf{s} above is hard in the case of neural networks as various correlations exist among the outputs of different neurons. (Khim & Loh, 2018) get around this by ignoring correlations and assuming that each path in the network can be independently optimized using

a different perturbation \mathbf{s} , which can often lead to vacuous results. For instance, consider a one-layer network, i.e. a linear combination of neurons. The network as a whole might be a large-margin classifier and hence robust on most inputs. However, if each neuron is a weakly correlated feature, then it can be independently attacked to induce a large loss. As a result, the approach of (Khim & Loh, 2018) greatly overestimates the adversarial loss in this case.

We avoid making such approximations. This requires a deeper analysis as in Sections 4 and 6, finally leading to adversarial shattering that could directly provide dimension-independent bounds. There have also been recent efforts to understand inference and generalization in adversarial settings for finite perturbation sets (Feige et al., 2015; Attias et al., 2019). Finally, the recent work of Wei & Ma (2019) provides generalization bounds for robust classification via studying a notion of layer-wise margin. This result is incomparable to our work as we aim to directly characterize the adversarial Rademacher complexity.

2. Preliminaries

We will denote vectors as lowercase bold letters (e.g., \mathbf{x}) and matrices as uppercase bold (e.g., \mathbf{M}). The all-ones vector is denoted by $\mathbf{1}$ and Hölder conjugates by a star (e.g., r^*). For a matrix \mathbf{M} , the (p, q) -group norm is defined as $\|\mathbf{M}\|_{p,q} = \|(\|\mathbf{M}_1\|_1, \dots, \|\mathbf{M}_d\|_p)\|_q$, where the \mathbf{M}_i s are the columns of \mathbf{M} .

We focus on binary classification over examples in \mathbb{R}^d and adversarial perturbations measured in ℓ_r -norm for $r \geq 1$. Given a loss function $\ell: \mathbb{R} \rightarrow [0, c]$, we define the loss of a hypothesis $f: \mathbb{R}^d \rightarrow \mathbb{R}$ on a pair $(\mathbf{x}, y) \in \mathbb{R}^d \times \{+1, -1\}$ as $\ell_f(\mathbf{x}, y) = \ell(yf(\mathbf{x}))$. As in the familiar classification setting, given a sample $\mathcal{S} = ((\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m))$ drawn i.i.d. from a distribution \mathcal{D} over $\mathbb{R}^d \times \{+1, -1\}$, we define the empirical risk and the expected risk of a hypothesis f as

$$R_{\mathcal{S}}(f) = \frac{1}{m} \sum_{i=1}^m \ell_f(\mathbf{x}_i, y_i) \quad R(f) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [\ell_f(\mathbf{x}, y)].$$

Let \mathcal{F} be a family of functions mapping from \mathbb{R}^d to \mathbb{R} . Then, the *empirical Rademacher complexity* of \mathcal{F} for a sample $\mathcal{S} = (\mathbf{x}_1, \dots, \mathbf{x}_m)$, is defined by

$$\widehat{\mathfrak{R}}_{\mathcal{S}}(\mathcal{F}) = \mathbb{E}_{\sigma} \left[\sup_{f \in \mathcal{F}} \frac{1}{m} \sum_{i=1}^m \sigma_i f(\mathbf{x}_i) \right], \quad (1)$$

where $\sigma = (\sigma_1, \dots, \sigma_m)$ is a vector of i.i.d. Rademacher variables, that is independent uniform random variables taking values in $\{-1, +1\}$. The *Rademacher complexity* of \mathcal{F} , $\mathfrak{R}_m(\mathcal{F})$, is defined as the expectation of this quantity: $\mathfrak{R}_m(\mathcal{F}) = \mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} [\widehat{\mathfrak{R}}_{\mathcal{S}}(\mathcal{F})]$, where \mathcal{D} is a distribution over the input space \mathbb{R}^d . The empirical Rademacher complexity

is a key data-dependent complexity measure. For a family of functions \mathcal{F} taking values in $[0, 1]$, the following learning guarantee holds: for any $\delta > 0$, with probability at least $1 - \delta$ over the draw of a sample $\mathcal{S} \sim \mathcal{D}^m$, the following inequality holds for all $f \in \mathcal{F}$ (Mohri et al., 2018):

$$\mathbb{E}_{x \sim \mathcal{D}} [f(x)] \leq \mathbb{E}_{x \sim \mathcal{S}} [f(x)] + 2\widehat{\mathfrak{R}}_{\mathcal{S}}(\mathcal{F}) + 3\sqrt{\frac{\log \frac{2}{\delta}}{2m}}, \quad (2)$$

where we denote by $\mathbb{E}_{x \sim \mathcal{S}} [f(x)]$ the empirical average of f , that is $\mathbb{E}_{x \sim \mathcal{S}} [f(x)] = \frac{1}{m} \sum_{i=1}^m f(x_i)$. A similar inequality holds for the average Rademacher complexity $\mathfrak{R}_m(\mathcal{F}_p) = \mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} [\widehat{\mathfrak{R}}_{\mathcal{S}}(\mathcal{F})]$:

$$\mathbb{E}_{x \sim \mathcal{D}} [f(x)] \leq \mathbb{E}_{x \sim \mathcal{S}} [f(x)] + 2\mathfrak{R}_m(\mathcal{F}) + \sqrt{\frac{\log \frac{1}{\delta}}{2m}}.$$

Furthermore, the Rademacher complexity of a hypothesis set also appears as a lower bound in generalization. As an example, for a symmetric family of functions \mathcal{F} taking values in $[-1, +1]$, the following holds (van der Vaart & Wellner, 1996):

$$\frac{1}{2} \left[\mathfrak{R}_m(\mathcal{F}) - \frac{1}{\sqrt{m}} \right] \leq \mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \sup_{f \in \mathcal{F}} \left| \mathbb{E}_{x \sim \mathcal{D}} [f(x)] - \mathbb{E}_{x \sim \mathcal{S}} [f(x)] \right| \leq 2\mathfrak{R}_m(\mathcal{F}).$$

An important application of these bounds is the derivation of margin bounds, which are crucial in the analysis of classification. Fix $\rho > 0$. Then, for any $\delta > 0$, with probability at least $1 - \delta$ over the draw of a sample $\mathcal{S} \sim \mathcal{D}^m$, the following inequality holds for all $f \in \mathcal{F}$ (Koltchinskii & Panchenko, 2002; Mohri et al., 2018):

$$\begin{aligned} & \mathbb{E}_{(x,y) \sim \mathcal{D}} [1_{yf(x) \leq 0}] \\ & \leq \frac{1}{m} \sum_{i=1}^m 1_{y_i f(x_i) \leq \rho} + \frac{2}{\rho} \widehat{\mathfrak{R}}_{\mathcal{S}}(\mathcal{F}) + 3\sqrt{\frac{\log \frac{2}{\delta}}{2m}}. \end{aligned} \quad (3)$$

Finer margin guarantees were recently presented by Cortes et al. (2020) in terms of Rademacher complexity and other complexity measures.

Robust Classification. We now extend the definitions above to their adversarial counterparts. In the setting of adversarially robust classification, the loss at (\mathbf{x}, y) is measured in terms of the worst loss incurred over an adversarial perturbation of \mathbf{x} within an ball of a certain radius in a norm $\|\cdot\|$. There are many possible norms we could use to measure a perturbation. A fairly general way to measure perturbations is the ℓ_r norm denoted $\|\cdot\|_r$. We will denote by ϵ the maximum magnitude of the allowed perturbations. Given $\epsilon > 0$, $r \geq 1$, a data point (\mathbf{x}, y) , a function $f: \mathbb{R}^d \rightarrow \mathbb{R}$, and a loss function $\ell: \mathbb{R} \rightarrow [0, c]$ we define the adversarial loss of f at (\mathbf{x}, y) as

$$\tilde{\ell}_f(\mathbf{x}, y) = \sup_{\|\mathbf{x} - \mathbf{x}'\|_r \leq \epsilon} \ell(yf(\mathbf{x}')).$$

Similarly, we define the adversarial empirical risk and the adversarial expected risk of a hypothesis f for a sample \mathcal{S} as follows:

$$\tilde{R}_{\mathcal{S}}(f) = \frac{1}{m} \sum_{i=1}^m \tilde{\ell}_f(\mathbf{x}_i, y_i) \quad \tilde{R}(f) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [\tilde{\ell}_f(\mathbf{x}, y)].$$

We also define $\tilde{\mathfrak{R}}(\mathcal{F})$, the *adversarial Rademacher complexity*, as the adversarial version of Rademacher complexity:

$$\tilde{\mathfrak{R}}_{\mathcal{S}}(\mathcal{F}) = \mathbb{E}_{\sigma} \left[\sup_{f \in \mathcal{F}} \frac{1}{m} \sum_{i=1}^m \sigma_i \sup_{\|\mathbf{x}_i - \mathbf{x}'_i\|_r \leq \epsilon} f(\mathbf{x}'_i) \right].$$

With the above definitions, the following is a consequence of (2) above and Talagrand's contraction lemma (Ledoux & Talagrand, 1991).

Theorem 1. *Let \mathcal{S} , \mathcal{D} , δ , and \mathcal{F} be as in equation (2). Further let ℓ be a Lipschitz loss function and define the class*

$$\ell_{\mathcal{F}} = \{\ell \circ f : f \in \mathcal{F}\}.$$

Then, with probability at least $1 - \delta$ the following holds for all $f \in \mathcal{F}$:

$$\tilde{R}(f) \leq \tilde{R}_{\mathcal{S}}(f) + 2c \cdot \tilde{\mathfrak{R}}_{\mathcal{S}}(\ell_{\mathcal{F}}) + 3c \sqrt{\frac{\log \frac{2}{\delta}}{2m}}. \quad (4)$$

Throughout the paper, we will assume that the loss function ℓ is non-increasing, a property satisfied by many common loss functions including the hinge loss, logistic loss and the exponential loss. In that case, as pointed out by Yin et al. (2019), the following equality holds:

$$\sup_{\|\mathbf{x}_i - \mathbf{x}'_i\|_r \leq \epsilon} \ell(y_i f(\mathbf{x}'_i)) = \ell \left(\inf_{\|\mathbf{x}_i - \mathbf{x}'_i\|_r \leq \epsilon} y_i f(\mathbf{x}'_i) \right).$$

Furthermore, when $\ell(\cdot)$ is L -Lipschitz, by Talagrand's Lemma, we have $\tilde{\mathfrak{R}}_{\mathcal{S}}(\ell_{\mathcal{F}}) \leq L \mathfrak{R}_{\mathcal{S}}(\tilde{\mathcal{F}})$, where $\tilde{\mathcal{F}}$ is the class defined by

$$\tilde{\mathcal{F}} = \{(\mathbf{x}, y) \mapsto \inf_{\|\mathbf{x} - \mathbf{x}'\|_r \leq \epsilon} y f(\mathbf{x}') : f \in \mathcal{F}\}. \quad (5)$$

Thus, we obtain the following inequalities:

$$\begin{aligned} \tilde{\mathfrak{R}}_{\mathcal{S}}(\ell_{\mathcal{F}}) &\leq L \mathfrak{R}_{\mathcal{S}}(\tilde{\mathcal{F}}) \\ &= L \mathbb{E}_{\sigma} \left[\sup_{f \in \mathcal{F}} \frac{1}{m} \sum_{i=1}^m \sigma_i \inf_{\|\mathbf{x}_i - \mathbf{x}'_i\|_r \leq \epsilon} y_i f(\mathbf{x}'_i) \right]. \end{aligned} \quad (6)$$

Providing sharp bounds on $\mathfrak{R}_{\mathcal{S}}(\tilde{\mathcal{F}})$ for various function classes \mathcal{F} will be the central focus of this work.

As an application, we can derive *robust margin bounds*:

Theorem 2 (Robust margin bounds). *Let $\mathcal{F}, \mathcal{S}, \mathcal{D}$, and δ be as in equation (3). Further let be $\tilde{\mathcal{F}}$ as in equation (5). Then, with probability at least $1 - \delta$ the following holds for all $f \in \mathcal{F}$:*

$$\tilde{R}(f) \leq \tilde{R}_{\mathcal{S}, \rho}(f) + \frac{2}{\rho} \mathfrak{R}_{\mathcal{S}}(\tilde{\mathcal{F}}) + 3 \sqrt{\frac{\log \frac{2}{\delta}}{2m}}. \quad (7)$$

3. Adversarial Rademacher Complexity of Linear Hypotheses

In this section, we provide a sharp characterization of the adversarial Rademacher complexity, as defined in (6), for linear function classes with bounded ℓ_p -norm and with perturbations measured in any ℓ_r -norm. Prior work (Yin et al., 2019) studied the case where the perturbations are measured in the ℓ_{∞} -norm. Our general analysis leads to a deeper understanding of the interplay between the complexity of the hypothesis classes (measured in ℓ_p -norm) and the perturbation set (measured in ℓ_r -norm), and how this dictates whether one can expect an additional dimension dependent penalty in the adversarial case over its non-adversarial counterpart. Furthermore, our analysis explicitly characterizes the dimension dependent term on which the adversarial Rademacher complexity depends on. This provides a finer analysis than the work of Khim & Loh (2018) and also has algorithmic implications. Formally, we study the case when

$$\mathcal{F}_p = \{\mathbf{x} \mapsto \langle \mathbf{w}, \mathbf{x} \rangle : \|\mathbf{w}\|_p \leq W\}. \quad (8)$$

3.1. Rademacher Complexity of Linear Hypotheses

A crucial aspect of our analysis in the linear case is a more general upper bound on the Rademacher complexity of p -norm bounded linear function classes, in the non-adversarial case. We first state this general bound as it will play an important role in later sections when analyzing the adversarial Rademacher complexity of ReLU functions and more general neural networks.

Theorem 3. *Let \mathcal{F}_p be the class of functions defined in (8). Then, given a sample $\mathcal{S} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)\}$ we have*

$$\mathfrak{R}_{\mathcal{S}}(\mathcal{F}_p) \leq \begin{cases} \frac{W}{m} \sqrt{2 \log(2d)} \|\mathbf{X}^{\top}\|_{2, p^*} & \text{if } p = 1 \\ \frac{\sqrt{2}W}{m} \left[\frac{\Gamma(\frac{p^*+1}{2})}{\sqrt{\pi}} \right]^{\frac{1}{p^*}} \|\mathbf{X}^{\top}\|_{2, p^*} & \text{if } 1 < p \leq 2 \\ \frac{W}{m} \|\mathbf{X}^{\top}\|_{2, p^*} & \text{if } p \geq 2 \end{cases}$$

Here, \mathbf{X} is the $d \times m$ matrix with the data points \mathbf{x}_i as columns. We make a few remarks about the theorem above and defer its proof to Appendix A.2. Some well-known bounds on the Rademacher complexity of \mathcal{F}_p are

$$\mathfrak{R}_{\mathcal{S}}(\mathcal{F}_p) \leq \begin{cases} W \sqrt{\frac{2 \log(2d)}{m}} \|\mathbf{X}\|_{\max} & \text{if } p = 1 \\ \frac{W}{m} \sqrt{p^* - 1} \|\mathbf{X}\|_{p^*, 2} & \text{if } 1 < p \leq 2 \end{cases} \quad (9)$$

Although the case $p \in [1, 2]$ in the theorem above is known (Kakade et al., 2008; Mohri et al., 2018), we provide a simpler proof of in Appendix A.1. The inequality for $p = 1$ is

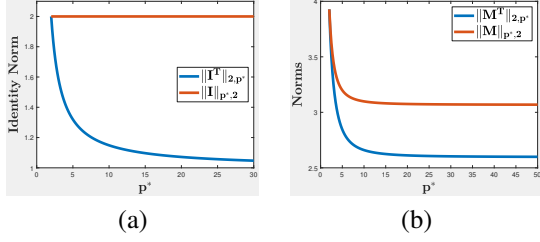


Figure 1: (a) A plot comparing two norms of the 4×4 identity matrix, $\|\mathbf{I}^\top\|_{2,p^*}$ and $\|\mathbf{I}\|_{p^*,2}$; the lower bound on the ratio of the two norms (10) in Proposition 1 holds for this matrix. (b) Same as (a), but for Gaussian matrices.

further reproduced for completeness. Our new bound coincides with (9) when $p = 2$ and is strictly better otherwise. Readers familiar with Rademacher complexity bounds for linear functions will notice that our bound in this case depends on the norm $\|\mathbf{X}^\top\|_{2,p^*}$. In contrast, standard bounds on the Rademacher complexity of linear classes depend on $\|\mathbf{X}\|_{p^*,2}$. In fact one can show that the $\|\mathbf{X}^\top\|_{2,p^*}$ is always smaller than $\|\mathbf{X}\|_{p^*,2}$ for $p \in (1, 2]$, that is $p^* \geq 2$, as shown by the last inequality of (10) in the following proposition.

Proposition 1. *Let \mathbf{M} be a $d \times m$ matrix. If $q \leq p$, then*

$$\min(m, d)^{\frac{1}{p} - \frac{1}{q}} \|\mathbf{M}^\top\|_{p,q} \leq \|\mathbf{M}\|_{q,p} \leq \|\mathbf{M}^\top\|_{p,q} \quad (10)$$

If $q \geq p$, then

$$\min(m, d)^{\frac{1}{p} - \frac{1}{q}} \|\mathbf{M}^\top\|_{p,q} \geq \|\mathbf{M}\|_{q,p} \geq \|\mathbf{M}^\top\|_{p,q} \quad (11)$$

These bounds are tight.

The proof is deferred to Appendix A.4. To visualize the ratio between these two norms, we plot the two norms for various values of p^* in Figure 1. For convenience, in the discussion below, we set $c_1(p) = \sqrt{p^* - 1}$ and $c_2(p) = \sqrt{2} \left[\frac{\Gamma(\frac{p^*+1}{2})}{\sqrt{\pi}} \right]^{\frac{1}{p^*}}$. Regarding the growth of the constant in our bound, one can show that as $p^* \rightarrow \infty$, $c_2(p)$ grows asymptotically like $e^{-\frac{1}{2}} \sqrt{p^*}$. In fact one can show that

$$e^{-\frac{1}{2}} \sqrt{p^*} \leq c_2(p) \leq e^{-\frac{1}{2}} \sqrt{p^* + 1}$$

Furthermore, $c_2(p) \leq c_1(p)$ in the relevant region (see Appendix A.5). In Figure 2 we plot $c_1(p)$, $c_2(p)$ and the bounds on $c_2(p)$ to illustrate the growth rate of these constants with p^* . Proposition 1 and that $c_2(p) \leq c_1(p)$ imply that the bounds we give for linear classes are stronger than what was previously known. We believe that the result in Theorem 3 is of wider interest and a recent companion note (Awasthi et al., 2020) provides a more detailed and self contained exposition.

3.2. Adversarial Rademacher Complexity of Linear Hypotheses

We now extend our bounds from the previous section to provide a complete characterization of the adversarial

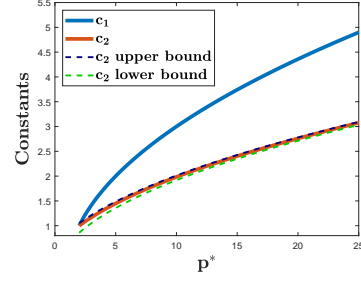


Figure 2: A plot of $c_1(p)$, $c_2(p)$, and the bounds from Lemma 4. Note that $c_1(2) = c_2(2)$ and that the upper and lower bounds on c_2 are tight.

Rademacher complexity of linear function classes under arbitrary r -norm perturbations. These theorems improve upon the recent work of Yin et al. (2019) that studies ∞ -norm perturbations and provide a finer analysis, with a matching lower bound, as compared to the recent work of Khim & Loh (2018). Our main result is stated below.

Theorem 4. *Let $\epsilon > 0$, $p, r \geq 1$. Consider a sample $\mathcal{S} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)\}$ with $\mathbf{x}_i \in \mathbb{R}^d$ and $y_i \in \{\pm 1\}$. Let \mathcal{F}_p be the class of linear functions defined in (8). Then it holds that*

$$\tilde{\mathfrak{R}}_{\mathcal{S}}(\mathcal{F}_p) \leq \left(\mathfrak{R}_{\mathcal{S}}(\mathcal{F}_p) + \epsilon \frac{W}{2\sqrt{m}} \max(d^{1-\frac{1}{r}-\frac{1}{p}}, 1) \right)$$

and

$$\tilde{\mathfrak{R}}_{\mathcal{S}}(\mathcal{F}_p) \geq \max \left(\mathfrak{R}_{\mathcal{S}}(\mathcal{F}_p), W \frac{\epsilon \max(d^{1-\frac{1}{r}-\frac{1}{p}}, 1)}{2\sqrt{2m}} \right).$$

Notice that when the perturbation is measured in ℓ_∞ -norm, i.e. $r = \infty$, we recover the bound of Yin et al. (2019). Hence the theorem above is a strict generalization of the result of Yin et al. (2019). Furthermore, when $\epsilon = 0$, as expected, the adversarial Rademacher complexity equals the standard Rademacher complexity of linear models and we can use our improved bounds from Theorem 3. The theorem above has important implications for the design of regularizers in the context of adversarial learning of linear models. As suggested by the upper bounds above, if $1/r + 1/p \geq 1$, then one can indeed perform adversarially robust learning with minimal statistical overhead in the standard classification setting! More specifically, in this case the upper bound on the adversarial Rademacher complexity has at most $W\epsilon/\sqrt{m}$ overhead on top of the standard bound from Theorem 3 and is dimension independent. Noting that $1 - 1/r = 1/r^*$, we get that for statistical efficiency one should choose a p -norm regularizer on \mathbf{w} , where $p \in [1, r^*]$. Our lower bound on the other hand shows that any other choice of a ℓ_p -norm based regularizer will necessarily incur a dimension-dependent penalty.

3.3. Proof sketch of Theorem 4

We provide a brief sketch of the proof of Theorem 4 and provide the details in Appendix B. As a first step, a simple argument shows that

$$\inf_{\|\mathbf{x}-\mathbf{x}'\|_r \leq \epsilon} y(\mathbf{w} \cdot \mathbf{x}') = y\mathbf{w} \cdot \mathbf{x} - \epsilon \|\mathbf{w}\|_{r^*}.$$

Using the above, we can write the adversarial Rademacher complexity as:

$$\mathfrak{R}_{\mathcal{S}}(\tilde{\mathcal{F}}_p) = \mathbb{E}_{\sigma} \left[\sup_{\|\mathbf{w}\|_p \leq W} \langle \mathbf{w}, \mathbf{u}_{\sigma} \rangle - \epsilon v_{\sigma} \|\mathbf{w}\|_{r^*} \right] \quad (12)$$

where, for convenience, we set $\mathbf{u}_{\sigma} = \frac{1}{m} \sum_{i=1}^m y_i \sigma_i \mathbf{x}_i$, $v_{\sigma} = \frac{1}{m} \sum_{i=1}^m \sigma_i$. Next, we present two key lemmas.

Lemma 1. *Let $1 \leq p, r \leq \infty$ and let d be the dimension. Then*

$$\sup_{\|\mathbf{w}\|_p \leq 1} \|\mathbf{w}\|_{r^*} = \max(1, d^{1-\frac{1}{r}-\frac{1}{p}})$$

Lemma 2. *Let $v_{\sigma} = \frac{1}{m} \sum_{i=1}^m \sigma_i$. Then it holds that*

$$\mathbb{E}_{\sigma} \left[\sup_{\|\mathbf{w}\|_p \leq W} v_{\sigma} \|\mathbf{w}\|_{r^*} \right] \geq \frac{W\epsilon \max(d^{1-\frac{1}{r}-\frac{1}{p}}, 1)}{2\sqrt{2m}},$$

and

$$\mathbb{E}_{\sigma} \left[\sup_{\|\mathbf{w}\|_p \leq W} v_{\sigma} \|\mathbf{w}\|_{r^*} \right] \leq \frac{W\epsilon \max(d^{1-\frac{1}{r}-\frac{1}{p}}, 1)}{2\sqrt{m}}$$

For the upper bound, using the sub-additivity of supremum and Lemma 2 yields

$$\begin{aligned} \mathfrak{R}_{\mathcal{S}}(\tilde{\mathcal{F}}_p) &\leq \mathfrak{R}_{\mathcal{S}}(\mathcal{F}_p) + \epsilon \mathbb{E}_{\sigma} \left[\sup_{\|\mathbf{w}\|_p \leq W} v_{\sigma} \|\mathbf{w}\|_{r^*} \right] \\ &= \mathfrak{R}_{\mathcal{S}}(\mathcal{F}_p) + \frac{1}{2} \epsilon \frac{W}{\sqrt{m}} \max(d^{1-\frac{1}{r}-\frac{1}{p}}, 1). \end{aligned}$$

For the lower bound, we apply two symmetrization arguments and show that

$$\mathfrak{R}_{\mathcal{S}}(\tilde{\mathcal{F}}_p) = \mathbb{E}_{\sigma} \left[\sup_{\|\mathbf{w}\|_p \leq W} -\langle \mathbf{w}, \mathbf{u}_{\sigma} \rangle + \epsilon v_{\sigma} \|\mathbf{w}\|_{r^*} \right] \quad (13)$$

$$= \mathbb{E}_{\sigma} \left[\sup_{\|\mathbf{w}\|_p \leq W} \langle \mathbf{w}, \mathbf{u}_{\sigma} \rangle + \epsilon v_{\sigma} \|\mathbf{w}\|_{r^*} \right]. \quad (14)$$

Averaging equations (12) and (14) and applying the sub-additivity of supremum gives:

$$\begin{aligned} \mathfrak{R}_{\mathcal{S}}(\tilde{\mathcal{F}}_p) &= \frac{1}{2} \mathbb{E}_{\sigma} \left[\sup_{\|\mathbf{w}\|_p \leq W} \langle \mathbf{w}, \mathbf{u}_{\sigma} \rangle - \epsilon v_{\sigma} \|\mathbf{w}\|_{r^*} \right] \\ &\quad + \frac{1}{2} \mathbb{E}_{\sigma} \left[\sup_{\|\mathbf{w}\|_p \leq W} \langle \mathbf{w}, \mathbf{u}_{\sigma} \rangle + \epsilon v_{\sigma} \|\mathbf{w}\|_{r^*} \right] \\ &\geq \mathbb{E}_{\sigma} \left[\sup_{\|\mathbf{w}\|_p \leq W} \langle \mathbf{w}, \mathbf{u}_{\sigma} \rangle \right] = W \mathfrak{R}_{\mathcal{S}}(\mathcal{F}_p). \end{aligned}$$

Now averaging (13) and (14), applying sub-additivity and Lemma 2, the following holds:

$$\begin{aligned} \mathfrak{R}_{\mathcal{S}}(\tilde{\mathcal{F}}_p) &= \frac{1}{2} \mathbb{E}_{\sigma} \left[\sup_{\|\mathbf{w}\|_p \leq W} -\langle \mathbf{w}, \mathbf{u}_{\sigma} \rangle + \epsilon v_{\sigma} \|\mathbf{w}\|_{r^*} \right] \\ &\quad + \frac{1}{2} \mathbb{E}_{\sigma} \left[\sup_{\|\mathbf{w}\|_p \leq W} \langle \mathbf{w}, \mathbf{u}_{\sigma} \rangle + \epsilon v_{\sigma} \|\mathbf{w}\|_{r^*} \right] \\ &\geq \mathbb{E}_{\sigma} \left[\sup_{\|\mathbf{w}\|_p \leq W} v_{\sigma} \|\mathbf{w}\|_{r^*} \right] \\ &\geq \frac{W}{2\sqrt{2m}} \epsilon \max(d^{1-\frac{1}{p}-\frac{1}{r}}, 1). \end{aligned}$$

4. Adversarial Rademacher Complexity of a Rectified Linear Unit

As a first step towards providing a bound for neural networks, in this section we study the adversarial Rademacher complexity of linear functions composed with a rectified linear unit (ReLU). Again we measure the size of functions in p norm and define the function class by

$$\mathcal{G}_p = \{(\mathbf{x}, y) \mapsto (y \langle \mathbf{w}, \mathbf{x} \rangle)_+ : \|\mathbf{w}\|_p \leq W, y \in \{-1, 1\}\} \quad (15)$$

where $(z)_+ = \max(z, 0)$. The following theorem presents a data-dependent upper bound on the adversarial Rademacher complexity of the ReLU unit.

Theorem 5. *Let \mathcal{G}_p be the class as defined in (15) and let \mathcal{F}_p be the corresponding linear class as defined in (8). Then, given a sample $\mathcal{S} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)\}$, the adversarial Rademacher complexity of \mathcal{G}_p can be bounded as follows:*

$$\tilde{\mathfrak{R}}_{\mathcal{S}}(\mathcal{G}_p) \leq \mathfrak{R}_{T_{\epsilon}}(\mathcal{F}_p) + \epsilon \frac{W}{2\sqrt{m}} \max(1, d^{1-\frac{1}{r}-\frac{1}{p}}),$$

where $T_{\epsilon} = \{i: y_i = -1 \text{ or } y_i = 1 \text{ and } \|\mathbf{x}_i\|_r > \epsilon\}$.

The second term in the bound above is similar to the dimension dependent term that appears in the linear case. The first term is the empirical Rademacher complexity of linear classes with bounded p -norm, but only measured on a carefully chosen subset of the data. This implies that data points with positive labels that have small norm as compared to the perturbation ϵ do not affect the Rademacher complexity. Hence, the guarantee in the theorem treats the two classes +1 and -1 asymmetrically.

This phenomenon originates from a property of the function $(z)_+$. Recall that in our setup $(z)_+$ will later be composed with a loss function ℓ . Because $\ell(yz_+)$ is the penalty incurred to the loss, the value yz_+ should be interpreted as a margin. Since the function $\max(0, z)$ is always 0 for $z \leq 0$, decreasing z below 0 does not affect the the margin. On the other hand increasing z above zero will increase the margin.

A large margin for a point labeled -1 corresponds to making z_+ as small as possible. As a result, every z with $z \leq 0$ gives the same margin. However, there is no upper bound on the margin for points in the class $+1$. As a result, the classifier $(y, \mathbf{x}) \mapsto y(\langle \mathbf{w}, \mathbf{x} \rangle)_+$ treats all non-negative margins for the class -1 in the same manner, but gives a higher reward for larger margins for the class $+1$.

This observation has implications for adversarial classification; as shown in Appendix C, an adversarially perturbed ReLU is $y \max(\mathbf{w} \cdot \mathbf{x} - \epsilon y \|\mathbf{w}\|_{r^*}, 0)$. If $\mathbf{w} \cdot \mathbf{x}$ is very negative, which corresponds to high confidence for $y = -1$, then the perturbation would not change the value of the loss function. On the other hand, if $\mathbf{w} \cdot \mathbf{x}$ were large and positive, a perturbation would definitely change the value of the margin and then influence the loss. We next complement our upper bound with a data dependent lower bound, stated below, on the adversarial Rademacher complexity.

Theorem 6. *Let \mathcal{G}_p be the class as defined in (15). Then it holds that*

$$\tilde{\mathfrak{R}}_{\mathcal{S}}(\mathcal{G}_p) \geq \frac{W}{2\sqrt{2}m} \sup_{\|\mathbf{s}\|_p=1} \left(\sum_{i \in T_{\epsilon, \mathbf{s}}} (\langle \mathbf{s}, \mathbf{x}_i \rangle - \epsilon y_i \|\mathbf{s}\|_{r^*})^2 \right)^{\frac{1}{2}}$$

where $T_{\epsilon, \mathbf{s}} = \{i: \langle \mathbf{s}, \mathbf{x}_i \rangle - y_i \epsilon \|\mathbf{s}\|_{r^*} > 0\}$.

A natural question that comes to mind is if one can characterize scenarios where the above lower bound leads to a dimension dependent term, as in the lower bound for linear hypotheses. In order to characterize this, for a given \mathbf{s} and $\delta > 0$, define the set $T_{\epsilon, \mathbf{s}}^{\delta}$ as

$$T_{\epsilon, \mathbf{s}}^{\delta} = \{i: \langle \mathbf{s}, \mathbf{x}_i \rangle - (1 + \delta y_i) y_i \epsilon \|\mathbf{s}\|_{r^*} > 0\}.$$

Notice that $T_{\epsilon, \mathbf{s}}^{\delta}$ is a subset of $T_{\epsilon, \mathbf{s}}$ and contains points in $T_{\epsilon, \mathbf{s}}$ that have a non-trivial margin. Then we get that

$$\begin{aligned} \mathfrak{R}_{\mathcal{S}}(\mathcal{G}_p) &\geq \frac{W}{2\sqrt{2}m} \sup_{\|\mathbf{s}\|_p=1} \left(\sum_{i \in T_{\epsilon, \mathbf{s}}^{\delta}} (\langle \mathbf{s}, \mathbf{x}_i \rangle - \epsilon y_i \|\mathbf{s}\|_{r^*})^2 \right)^{\frac{1}{2}} \\ &\geq \frac{W}{2\sqrt{2}m} \sup_{\|\mathbf{s}\|_p=1} \left(\sum_{i \in T_{\epsilon, \mathbf{s}}^{\delta}} (\delta \epsilon \|\mathbf{s}\|_{r^*})^2 \right)^{\frac{1}{2}} \\ &= \frac{W \delta \epsilon}{2\sqrt{2}m} \sup_{\|\mathbf{s}\|_p=1} |T_{\epsilon, \mathbf{s}}^{\delta}| \|\mathbf{s}\|_{r^*}. \end{aligned}$$

Denoting \mathbf{s}^* to be the vector that achieves the value $\sup_{\|\mathbf{s}\|_p=1} \|\mathbf{s}\|_{r^*}$ we get that

$$\mathfrak{R}_{\mathcal{S}}(\mathcal{G}_p) \geq \frac{W \delta \epsilon}{2\sqrt{2}m} |T_{\epsilon, \mathbf{s}^*}^{\delta}| \max(d^{1-\frac{1}{p}-\frac{1}{r}}, 1).$$

Hence, if for a given constant $\delta > 0$, the size of the set $T_{\epsilon, \mathbf{s}^*}^{\delta}$ is large then we expect a dimension dependent lower bound similar to the linear case.

5. Adversarial Rademacher Complexity of Neural Nets

Building on our analysis for the case of a single ReLU unit, we next give an upper bound on the adversarial Rademacher complexity for the class of one-layer neural networks comprised of a Lipschitz activation ρ with $\rho(0) = 0$. The guarantees of our theorem resemble the bound on the standard Rademacher complexity of neural networks, as provided in (Cortes et al., 2017). An analysis based on other forms of generalization bounds on neural nets is also possible, such as that of Bartlett et al. (2017). The family of functions of such one-layer neural networks is defined as follows:

$$\mathcal{G}_p^n = \left\{ (\mathbf{x}, y) \mapsto y \sum_{j=1}^n u_j \rho(\mathbf{w}_j \cdot \mathbf{x}) : \|\mathbf{u}\|_1 \leq \Lambda, \|\mathbf{w}_j\|_p \leq W \right\}.$$

Our main theorem is stated below.

Theorem 7. *Let ρ be a function with Lipschitz constant L_{ρ} with $\rho(0) = 0$ and consider perturbations in r -norm. Then, the following upper bound holds for the adversarial Rademacher complexity of \mathcal{G}_p^n :*

$$\tilde{\mathfrak{R}}_{\mathcal{S}}(\mathcal{G}_p^n) \leq L_{\rho} \left[\frac{W \Lambda \max(1, d^{1-\frac{1}{p}-\frac{1}{r}}) (\|\mathbf{X}\|_{r, \infty} + \epsilon)}{\sqrt{m}} \right] \times \left(1 + \sqrt{d(n+1) \log(36)} \right).$$

The proof is presented in Appendix D. The only requirements on our activation function ρ is that it is Lipschitz and $\rho(0) = 0$. This stipulation is satisfied by common activation functions like the ReLU, the leaky ReLU, and the hyperbolic tangent, but not the sigmoid or a step function. In comparison to the adversarial Rademacher complexity of linear classifiers, Theorem 7 still includes a $\max(1, d^{1-\frac{1}{p}-\frac{1}{r}})$ factor, again implying that one should choose a model class with $p \leq r^*$. The complexity of the vector \mathbf{u} is bounded by ℓ_1 norm as that is what turns out to be natural in the proof. However, the dimension dependence is larger by a factor of \sqrt{d} . The dependence on the number of neurons (\sqrt{n}) is also problematic. This fact is unfortunate since a much larger sample size m would be required for good generalization. In the next section we present a promising approach towards removing the dependence on dimension and the number of neurons in the above bound.

6. Towards Dimension-Independent Bounds

In this section we introduce a new framework for analyzing the adversarial Rademacher complexity of neural networks with ReLU activations. Unlike the case of linear hypotheses, the dimension-dependent term in the upper bound in Theorem 7 cannot be avoided by simply picking the appropriate norm p . In particular, deriving dimension-independent

bounds for the adversarial Rademacher complexity of neural networks is a difficult problem. Prior works (Yin et al., 2019; Khim & Loh, 2018) have resorted to bounding the adversarial Rademacher complexity of surrogates that are more tractable. However, it is not clear how those guarantees translate into meaningful bounds on the complexity of the original network. In this section, we present an approach towards obtaining dimension-independent bounds on the adversarial Rademacher complexity of the original network.

A major component of the difficulty in analyzing adversarial Rademacher complexity relates to providing a tight characterization of the optimal adversarial perturbation for a given point \mathbf{x}_i , i.e.,

$$\mathbf{s}_i^* = \operatorname{argmin}_{\mathbf{s}: \|\mathbf{s}\|_r \leq 1} y_i \sum_{j=1}^n u_j (\mathbf{w}_j \cdot (\mathbf{x}_i + \epsilon \mathbf{s}))_+ \quad (16)$$

Thus, to begin, we study properties of such adversarial perturbations to the neural network. Afterwards, we leverage these properties to bound the adversarial Rademacher complexity. Notably, the proofs of these properties heavily rely on the fact that the activation function is ReLU and not any other Lipschitz function. As in the previous section, we will focus on the family of a one layer-network \mathcal{G}_p^n with activation $\rho(z) = z_+$.

6.1. Characterizing Adversarial Perturbations

In this section, we discuss characteristics of adversarial perturbations to neural networks with ReLU activations. The following theorem implies that, if the perturbations are bounded in ℓ_r -norm by ϵ , then typically the optimal adversarial perturbations will have exactly r -norm ϵ .

Theorem 8. *Let d be the dimension and n the number of neurons. Consider the problem*

$$\inf_{\|\mathbf{s}\|_r \leq 1} f(\mathbf{s}) = \sum_{j=1}^n u_j (\mathbf{w}_j \cdot (\mathbf{x} + \epsilon \mathbf{s}))_+. \quad (17)$$

If either $\|\mathbf{x}\|_r \geq \epsilon$ or $n < d$, an optimum is attained on the sphere $\{\mathbf{s}: \|\mathbf{s}\|_r = 1\}$. Otherwise, an optimum is attained either at $\mathbf{s} = -\frac{1}{\epsilon} \mathbf{x}$ or on $\|\mathbf{s}\|_r = 1$.

The proof of the above theorem is deferred to Appendix E.1. Theorem 8 implies that if $n < d$, then the optimal perturbation always has norm ϵ . This result is significant because $n < d$ is a common scenario. At the same time, the theorem also implies that if $\|\mathbf{x}\|_r \geq \epsilon$, then the optimal perturbation still has norm ϵ . In practice, one expects the norm of the data points to be larger than the perturbation. Thus, on real world datasets, one would expect adversarial perturbations to always have norm ϵ .

For $1 < r < \infty$, Theorem 8 aids in finding a necessary condition for the optimum. This condition implies that

critical points are characterized by specifying which \mathbf{w}_j satisfy $\mathbf{w}_j \cdot (\mathbf{x} + \epsilon \mathbf{s}) < 0$, $\mathbf{w}_j \cdot (\mathbf{x} + \epsilon \mathbf{s}) = 0$, and $\mathbf{w}_j \cdot (\mathbf{x} + \epsilon \mathbf{s}) > 0$. The exact assertion is fairly involved, so we delay the statement of this theorem to Appendix E.2. However, the theorem simplifies considerably for $r = 2$ and we include this case below.

Theorem 9. *Assume that $\|\mathbf{x}\|_r \geq \epsilon$. Let $1 < r < \infty$ and take f as in Theorem 8 and \mathbf{s}^* as the minimizer of (17). Define the following three sets:*

$$\begin{aligned} N &= \{j: \mathbf{w}_j \cdot (\mathbf{x} + \epsilon \mathbf{s}^*) < 0\} \\ Z &= \{j: \mathbf{w}_j \cdot (\mathbf{x} + \epsilon \mathbf{s}^*) = 0\} \\ P &= \{j: \mathbf{w}_j \cdot (\mathbf{x} + \epsilon \mathbf{s}^*) > 0\}. \end{aligned}$$

\mathbf{s}^ is characterized by specifying the sets N, Z , and P . Furthermore, if $r = 2$, \mathbf{s}^* can be explicitly expressed in terms of these sets. Let P_Z be the projection onto $\operatorname{span}\{\mathbf{w}_j\}_{j \in Z}$ and P_{Z^c} the projection onto the complement of this subspace. Then, \mathbf{s}^* is given by*

$$\mathbf{s}^* = - \left(\sqrt{1 - \frac{\|P_Z \mathbf{x}\|_2^2}{\epsilon^2}} \frac{P_{Z^c} \sum_{j \in P} u_j \mathbf{w}_j}{\|P_{Z^c} \sum_{j \in P} u_j \mathbf{w}_j\|_2} + \frac{1}{\epsilon} P_Z \mathbf{x} \right).$$

6.2. Dimension-Independent Bound for ReLU Neural Networks

Observe that, given \mathbf{u} and the weight matrix \mathbf{W} with columns $(\mathbf{w}_1, \dots, \mathbf{w}_n)$, each \mathbf{x}_i partitions these vectors into three sets depending on whether at the optimal \mathbf{s}_i^* , $\mathbf{w}_j \cdot (\mathbf{x}_i + \epsilon \mathbf{s}_i^*)$ is positive, zero or negative. As a result, given \mathbf{W} and \mathbf{u} , the points in the data set can be partitioned into sets depending on whether they induce the same sign pattern on the columns of \mathbf{W} . Let \mathcal{C}_S denote the set of all such possible partitions and let C_S^* be the size of this set. Indexing a particular partition in this set by \mathcal{C} , let $n_{\mathcal{C}}$ be the number of parts in this partition and define $\Pi_S^* = \max_{\mathcal{C}} n_{\mathcal{C}}$. Notice that both Π_S^* and C_S^* are data-dependent quantities. We next state a general theorem that does not explicitly depend on the dimension and instead bounds the adversarial Rademacher complexity in terms of the above data-dependent quantities.

Theorem 10. *Consider the family of functions \mathcal{G}_p^n with activation function $\rho(z) = (z)_+$ and perturbations in r -norm for $1 < r < \infty$. Assume that for our sample $\|\mathbf{x}_i\|_r \geq \epsilon$. Then, the following upper bound on the Rademacher complexity holds:*

$$\tilde{\mathfrak{R}}_S(\mathcal{G}_p^n) \leq \left[\frac{W \Lambda \max(1, d^{1-\frac{1}{p}-\frac{1}{r}})}{\sqrt{m}} (\|\mathbf{X}\|_{p^*, \infty} + \epsilon) \right] C_S^* \sqrt{\Pi_S^*}.$$

Notice that the main difference between the above guarantee and the one from the previous section is that the dimension-dependent term $(1 + \sqrt{d(n+1)} \log(9m))$ has been replaced by data-dependent quantities. Next, we discuss how to

bound these data-dependent quantities in terms of a notion of *adversarial shattering* that we introduce in this work.

Bounding $\Pi_{\mathcal{S}}^*$ and ϵ -adversarial shattering. A key quantity of interest in understanding the bounds from the above theorem is $\Pi_{\mathcal{S}}^*$. Notice that this corresponds to the maximum number of partitions of the vectors $\mathbf{w}_1, \dots, \mathbf{w}_j$ that can be induced by the dataset $(\mathbf{x}_1, \dots, \mathbf{x}_m)$. Viewing the \mathbf{w}_j s as examples and the \mathbf{x}_i s as hyperplanes, this corresponds to the number of sign patterns on \mathbf{W} that can be induced by \mathcal{S} . In standard settings, this would be bounded by the VC-dimension (d in this case). However, we know more about how the \mathbf{x}_i s act on these vectors. Notice that at the optimal \mathbf{s}_i^* for a given \mathbf{x}_i , for some subset of vectors $\mathbf{w}_j \cdot \mathbf{x}_i + \mathbf{w}_j \cdot \epsilon \mathbf{s}_i^* \geq 0$, and for the rest it must be that $\mathbf{w}_j \cdot \mathbf{x}_i + \mathbf{w}_j \cdot \epsilon \mathbf{s}_i^* \leq 0$. Hence, not only does \mathbf{x}_i induce a sign pattern on the \mathbf{w}_j s, it does so with a certain margin. This is reminiscent of the classical notion of *fat shattering* (Mohri et al., 2018) from statistical learning theory. However, in this case, the margin induced could itself depend on the \mathbf{w}_j s in a complex manner via the product of $\mathbf{w}_j \cdot \mathbf{s}_i^*$. To formalize this intuition, we define the following notion of ϵ -adversarial shattering.

Definition 1. Fix the sample $\mathcal{S} = ((\mathbf{x}_1, y_1) \dots (\mathbf{x}_m, y_m))$ and $(\mathbf{w}_1, \dots, \mathbf{w}_n)$. Let $\mathbf{s}_i = \operatorname{argmin}_{\|\mathbf{s}\|_r \leq 1} y_i \sum_{j=1}^n u_j(\mathbf{w}_j \cdot (\mathbf{x}_i + \epsilon \mathbf{s}))_+$, and define the following three sets:

$$\begin{aligned} P_i &= \{j: \mathbf{w}_j \cdot (\mathbf{x} + \epsilon \mathbf{s}_i) > 0\} \\ Z_i &= \{j: \mathbf{w}_j \cdot (\mathbf{x} + \epsilon \mathbf{s}_i) = 0\} \\ N_i &= \{j: \mathbf{w}_j \cdot (\mathbf{x} + \epsilon \mathbf{s}_i) < 0\}. \end{aligned}$$

Let $\Pi_{\mathcal{S}}(\mathbf{W})$ be the number of distinct (P_i, Z_i, N_i) s that are induced by \mathcal{S} , where \mathbf{W} is a matrix that admits the \mathbf{w}_j s as columns. We call $\Pi_{\mathcal{S}}(\mathbf{W})$ the ϵ -adversarial growth function. We say that \mathbf{W} is ϵ -adversarially shattered if every $P \subset [n]$ is possible.

Under certain assumptions, by carefully studying the above notion of adversarial shattering one can obtain bounds of the form $O(\frac{1}{\epsilon^2})$ on the maximum number of \mathbf{w}_j s that can be adversarially shattered by \mathcal{S} . This lets us use an argument similar in spirit to Sauer’s lemma (Sauer, 1972; Shelah, 1972) to bound $\Pi_{\mathcal{S}}^*$ by $n^{O(1/\epsilon^2)}$, thereby leading to a meaningful bound in Theorem 10. We believe that a further study of the above notion of adversarial shattering is the key to proving general dimension-independent bounds on the adversarial Rademacher complexity of neural networks.

7. Conclusion

In this work we presented a detailed study of the generalization properties of linear models and neural networks under adversarial perturbations. Our bounds for the linear case improve upon prior work and also lead to a novel analysis of the Rademacher complexity of linear hypotheses in

non-adversarial settings as well. For the case of a single ReLU unit, while we have upper and lower bounds, it would be interesting to investigate the extent to which they are close to each other. Our analysis for the linear and ReLU hypotheses reveals that by choosing the appropriate norm regularization (ℓ_p) on the weight matrices, one can indeed avoid dimension dependence and achieve generalization in adversarial settings with negligible statistical overhead as compared to the corresponding non-adversarial setting. Our analysis illustrates the importance of choosing p satisfying $\frac{1}{r} + \frac{1}{p} \geq 1$ in algorithms. This relationship further suggests that for robustness to perturbations in an arbitrary norm $\|\cdot\|$, one could regularize by the dual norm of $\|\cdot\|$. Investigating this relationship could be future work. Finally, it would be interesting to use our approach from Section 6.2 based on ϵ -adversarial shattering to provide dimension-independent upper bounds on the adversarial Rademacher complexity of neural networks.

References

- Alzer, H. On some inequalities for the Gamma and Psi functions. *Math. Comput.*, 66(217):373–389, 1997.
- Athalye, A., Carlini, N., and Wagner, D. Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. *arXiv preprint arXiv:1802.00420*, 2018.
- Attias, I., Kontorovich, A., and Mansour, Y. Improved generalization bounds for robust learning. In *Algorithmic Learning Theory*, pp. 162–183, 2019.
- Awasthi, P., Dutta, A., and Vijayaraghavan, A. On robustness to adversarial examples and polynomial optimization. In *NeurIPS*, pp. 13737–13747, 2019.
- Awasthi, P., Frank, N., and Mohri, M. On the Rademacher complexity of linear hypothesis sets. *CoRR*, abs/2007.11045, 2020. URL <https://arxiv.org/abs/2007.11045>.
- Bartlett, P. L., Foster, D. J., and Telgarsky, M. Spectrally-normalized margin bounds for neural networks. *CoRR*, 2017.
- Bubeck, S., Lee, Y. T., Price, E., and Razenshteyn, I. Adversarial examples from cryptographic pseudo-random generators. *arXiv preprint arXiv:1811.06418*, 2018a.
- Bubeck, S., Price, E., and Razenshteyn, I. Adversarial examples from computational constraints. *arXiv preprint arXiv:1805.10204*, 2018b.
- Carlini, N. and Wagner, D. Towards evaluating the robustness of neural networks. In *2017 IEEE Symposium on Security and Privacy (SP)*, pp. 39–57. IEEE, 2017.

- Carlini, N. and Wagner, D. Audio adversarial examples: Targeted attacks on speech-to-text. In *2018 IEEE Security and Privacy Workshops (SPW)*, pp. 1–7. IEEE, 2018.
- Chen, X., Liu, C., Li, B., Lu, K., and Song, D. Targeted backdoor attacks on deep learning systems using data poisoning. *arXiv preprint arXiv:1712.05526*, 2017.
- Cortes, C., Gonzalvo, X., Kuznetsov, V., Mohri, M., and Yang, S. AdaNet: Adaptive structural learning of artificial neural networks. In *Proceedings of ICML*, pp. 874–883, 2017.
- Cortes, C., Mohri, M., and Suresh, A. T. Relative deviation margin bounds. *CoRR*, abs/2006.14950, 2020.
- Degwekar, A., Nakkiran, P., and Vaikuntanathan, V. Computational limitations in robust classification and win-win results. In *Proceedings of the Thirty-Second Conference on Learning Theory*, pp. 994–1028, 2019.
- Eykholt, K., Evtimov, I., Fernandes, E., Li, B., Rahmati, A., Xiao, C., Prakash, A., Kohno, T., and Song, D. Robust physical-world attacks on deep learning visual classification. In *Proceedings of CVPR*, pp. 1625–1634, 2018.
- Feige, U., Mansour, Y., and Schapire, R. Learning and inference in the presence of corrupted inputs. In *Conference on Learning Theory*, pp. 637–657, 2015.
- Goodfellow, I. J., Shlens, J., and Szegedy, C. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014.
- Gowal, S., Dvijotham, K., Stanforth, R., Bunel, R., Qin, C., Uesato, J., Arandjelovic, R., Mann, T., and Kohli, P. On the effectiveness of interval bound propagation for training verifiably robust models. *arXiv preprint arXiv:1810.12715*, 2018.
- Gowal, S., Uesato, J., Qin, C., Huang, P.-S., Mann, T., and Kohli, P. An alternative surrogate loss for PGD-based adversarial testing. *arXiv preprint arXiv:1910.09338*, 2019.
- Haagerup, U. The best constants in the Khintchine inequality. *Studia Mathematica*, 70:231–283, 1981.
- Huber, P. J. *Robust statistics*. Springer, 2011.
- Kakade, S. M., Sridharan, K., and Tewari, A. On the complexity of linear prediction: Risk bounds, margin bounds, and regularization. In *Proceedings of NIPS*, pp. 793–800, 2008.
- Kearns, M. and Li, M. Learning in the presence of malicious errors. *SIAM Journal on Computing*, 22(4):807–837, 1993.
- Kearns, M. J., Schapire, R. E., and Sellie, L. M. Toward efficient agnostic learning. *Machine Learning*, 17(2-3): 115–141, 1994.
- 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. *Annals of Statistics*, 30, 2002.
- Ledoux, M. and Talagrand, M. *Probability in Banach Spaces: Isoperimetry and Processes*. Springer, New York, 1991.
- Madry, A., Makelov, A., Schmidt, L., Tsipras, D., and Vladu, A. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017.
- Mohri, M., Rostamizadeh, A., and Talwalkar, A. *Foundations of Machine Learning*. The MIT Press, second edition, 2018.
- Montasser, O., Hanneke, S., and Srebro, N. VC classes are adversarially robustly learnable, but only improperly. *arXiv preprint arXiv:1902.04217*, 2019.
- Nakkiran, P. Adversarial robustness may be at odds with simplicity. *arXiv preprint arXiv:1901.00532*, 2019.
- Olver, F. W. J., Lozier, D. W., Boisvert, R. F., and Clark, C. W. *The NIST Handbook of Mathematical Functions*. Cambridge Univ. Press, 2010.
- Polyakova, L. On minimizing the sum of a convex function and a concave function. *Mathematical Programming Study*, 29, 1986.
- Raghunathan, A., Xie, S. M., Yang, F., Duchi, J. C., and Liang, P. Adversarial training can hurt generalization. *arXiv preprint arXiv:1906.06032*, 2019.
- Sauer, N. On the density of families of sets. *Journal of Combinatorial Theory, Series A*, 13(1):145–147, 1972.
- Schmidt, L., Santurkar, S., Tsipras, D., Talwar, K., and Madry, A. Adversarially robust generalization requires more data. *arXiv preprint arXiv:1804.11285*, 2018.
- Schott, L., Rauber, J., Bethge, M., and Brendel, W. Towards the first adversarially robust neural network model on mnist. *arXiv preprint arXiv:1805.09190*, 2018.
- Sharma, Y. and Chen, P.-Y. Breaking the madry defense model with l_1 -based adversarial examples. *arXiv preprint arXiv:1710.10733*, 2017.

Shelah, S. A combinatorial problem; stability and order for models and theories in infinitary languages. *Pacific Journal of Mathematics*, 41(1):247–261, 1972.

Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I. J., and Fergus, R. Intriguing properties of neural networks. In *Proceedings of ICLR*, 2014.

Tsipras, D., Santurkar, S., Engstrom, L., Turner, A., and Madry, A. Robustness may be at odds with accuracy, 2018.

van der Vaart, A. W. and Wellner, J. A. *Weak Convergence and Empirical Processes*. Springer, 1996.

Wei, C. and Ma, T. Improved sample complexities for deep networks and robust classification via an all-layer margin. *arXiv preprint arXiv:1910.04284*, 2019.

Yin, D., Ramchandran, K., and Bartlett, P. L. Rademacher complexity for adversarially robust generalization. In *Proceedings of ICML*, pp. 7085–7094, 2019.