

---

# Why Robust Generalization in Deep Learning is Difficult: Perspective of Expressive Power

---

Binghui Li<sup>1,6,7,\*</sup> Jikai Jin<sup>2,\*</sup> Han Zhong<sup>3</sup> John E. Hopcroft<sup>4</sup> Liwei Wang<sup>3,5,†</sup>

<sup>1</sup>School of EECS, Peking University

<sup>2</sup>School of Mathematical Sciences, Peking University

<sup>3</sup>Center for Data Science, Peking University <sup>4</sup>Cornell University

<sup>5</sup>National Key Laboratory of General Artificial Intelligence,  
School of Intelligence Science and Technology, Peking University

<sup>6</sup>Peng Cheng Laboratory <sup>7</sup>Pazhou Laboratory (Huangpu)  
{libinghui, jkjin}@pku.edu.cn, hanzhong@stu.pku.edu.cn,  
jeh17@cornell.edu, wanglw@cis.pku.edu.cn

## Abstract

It is well-known that modern neural networks are vulnerable to adversarial examples. To mitigate this problem, a series of robust learning algorithms have been proposed. However, although the robust training error can be near zero via some methods, all existing algorithms lead to a high robust generalization error. In this paper, we provide a theoretical understanding of this puzzling phenomenon from the perspective of expressive power for deep neural networks. Specifically, for binary classification problems with well-separated data, we show that, for ReLU networks, while mild over-parameterization is sufficient for high robust training accuracy, there exists a constant robust generalization gap unless the size of the neural network is exponential in the data dimension  $d$ . This result holds even if the data is linear separable (which means achieving standard generalization is easy), and more generally for any parameterized function classes as long as their VC dimension is at most polynomial in the number of parameters. Moreover, we establish an improved upper bound of  $\exp(\mathcal{O}(k))$  for the network size to achieve low robust generalization error when the data lies on a manifold with intrinsic dimension  $k$  ( $k \ll d$ ). Nonetheless, we also have a lower bound that grows exponentially with respect to  $k$  — the curse of dimensionality is inevitable. By demonstrating an exponential separation between the network size for achieving low robust training and generalization error, our results reveal that the hardness of robust generalization may stem from the expressive power of practical models.

## 1 Introduction

Deep neural networks have achieved remarkable success in a variety of disciplines including computer vision (Voulodimos et al., 2018), natural language processing (Devlin et al., 2018) as well as scientific and engineering applications (Jumper et al., 2021). However, it is observed that neural networks are often sensitive to small adversarial attacks (Biggio et al., 2013; Szegedy et al., 2013; Goodfellow et al., 2014), which potentially gives rise to reliability and security problems in real-world applications.

In light of this pitfall, it is highly desirable to obtain classifiers that are robust to small but adversarial perturbations. A common approach is to design adversarial training algorithms by using adversarial

---

\*Equal contribution.

†Corresponding author.

examples as training data (Madry et al., 2017; Tramèr et al., 2018; Shafahi et al., 2019). Another line of works (Cohen et al., 2019; Zhang et al., 2021a) proposes some provably robust models to tackle this problem. However, while the state-of-the-art adversarial training methods can achieve high robust training accuracy (e.g. nearly 100% on CIFAR-10 (Raghunathan et al., 2019)), all existing methods suffer from large robust test error. Therefore, it is natural to ask what is the cause for such a large generalization gap in the context of robust learning.

Previous works have studied the hardness of achieving adversarial robustness from different perspectives. A well-known phenomenon called the *robustness-accuracy tradeoff* has been empirically observed (Raghunathan et al., 2019) and theoretically proven to occur in different settings (Tsipras et al., 2019; Zhang et al., 2019). Dohmatob (2019) shows that adversarial robustness is impossible to achieve under certain assumptions on the data distribution, while it is shown in Nakkiran (2019) that even when an adversarial robust classifier does exist, it can be exponentially more complex than its non-robust counterpart. Hassani and Javanmard (2022) studies the role of over-parameterization on adversarial robustness by focusing on random features regression models.

At first glance, these works seem to provide convincing evidence that robustness is hard to achieve in general. However, this view is challenged by Yang et al. (2020), who observes that for real data sets, different classes tend to be well-separated (as defined below), while the perturbation radius is often much smaller than the separation distance. As pointed out by Yang et al. (2020), all aforementioned works fail to take this separability property of data into consideration.

**Definition 1.1** (Separated data). *Suppose that  $A, B \subset \mathbb{R}^d$  and  $\epsilon > 0$ . We say that  $A, B$  are  $\epsilon$ -separated under  $\ell_p$  norm ( $1 \leq p \leq +\infty$ ) if*

$$\|\mathbf{x}_A - \mathbf{x}_B\|_p \geq \epsilon, \quad \forall \mathbf{x}_A \in A, \mathbf{x}_B \in B.$$

Indeed, this assumption is necessary to ensure the existence of a robust classifier. Without this separated condition, it is clear that there is no robust classifier even if a non-robust classifier always exists, as discussed above.

Recently, Bubeck and Sellke (2021) shows that for regression problems, over-parameterization may be necessary for achieving robustness. However, they measure robustness of a model by its training error and Lipschitz constant, which has a subtle difference with *robust test error* (Madry et al., 2017); see the discussions in (Bubeck and Sellke, 2021, Section 1.1).

To sum up the above, although existing robust training algorithms result in low robust test accuracy, previous works do not provide a satisfactory explanation of this phenomenon, since there exists a gap between their settings and practice. In particular, it is not known whether achieving robustness can be easier for data with additional structural properties such as separability (Yang et al., 2020) and low intrinsic dimensionality (Gong et al., 2019).

In this paper, we make an important step towards understanding robust generalization from the viewpoint of neural network expressivity. Focusing on binary classification problems with separated data (cf. Definition 1.1) in  $\mathbb{R}^d$ , we make the following contributions:

- Given a data set  $\mathcal{D}$  of size  $N$  that satisfies a separability condition, we show in Section 2 that it is possible for a ReLU network with  $\tilde{\mathcal{O}}(Nd)$  weights to robustly classify  $\mathcal{D}$ . In other words, an over-parameterized ReLU network with reasonable size can achieve 100% robust training accuracy.
- We next consider the robust test error (cf. Definition 3.1). As a warm-up, we show in Section 3 that, in contrast with the robust *training* error, mere separability of data does not imply that low robust test error can be attained by neural networks, unless their size is exponential in  $d$ . This motivates the subsequent sections where we consider data with additional structures.
- In Section 4, we prove the main result of this paper, which states that for achieving low robust test error, an  $\exp(\Omega(d))$  lower bound on the network size is inevitable, even when the underlying distributions of the two classes are linear separable. Moreover, this lower bound holds for arbitrarily small perturbation radius and more general models as long as their VC dimension is at most polynomial in the number of parameters.
- Finally, in Section 5 we consider data that lies on a  $k$ -dimensional manifold ( $k \ll d$ ), and prove an improved upper bound  $\exp(\mathcal{O}(k))$  for the size of neural networks for achieving

Table 1: Summary of our main results.

Params	Setting			
	Robust Training	Robust Generalization		
		General Case	Linear Separable	$k$ -dim Manifold
Upper Bound	$\mathcal{O}(Nd)$ (Thm 2.2)	$\exp(\mathcal{O}(d))$ (Thm 3.3)	$\exp(\mathcal{O}(k))$ (Thm 5.5)	
Lower Bound	$\Omega(\sqrt{Nd})$ (Thm 2.3)	$\exp(\Omega(d))$ (Thm 3.4)	$\exp(\Omega(d))$ (Thm 4.3)	$\exp(\Omega(k))$ (Thm 5.8)

low robust test error. Nonetheless, the curse of dimensionality is inescapable – the lower bound is also exponential in  $k$ .

The upper and lower bounds on network size are summarized in Table 1. Overall, our theoretical analysis suggests that the hardness of achieving robust generalization may stem from the expressive power of practical models.

### 1.1 Implications of our results

Before moving on to technical parts, we would like to first discuss the implications of our results by comparing them to previous works.

Our main result is the exponential lower bound on the neural network size for generalization. First, different from previous hardness results, our result is established for data sets that have desirable structural properties, hence more closely related to practical settings. Note that the separability condition implies that there *exists* a classifier that can perfectly and robustly classify the data i.e. achieve zero robust test error. However, we show that such classifier is hard to approximate using neural networks with moderate size.

Second, it is a popular belief that many real-world data sets are intrinsically low dimensional, although they lie in a high dimensional space. Our results imply that low dimensional structure makes robust generalization possible with a neural network with significantly smaller size (when  $k \ll d$ ). However, the size must still be exponential in  $k$ .

Finally, we show that there exists an *exponential* separation between the required size of neural networks for achieving low robust training and test error. Based on our results, we conjecture that the widely observed drop of robust test accuracy is not due to limitations of existing algorithms – rather, it is a more fundamental issue originating from the expressive power of neural networks.

### 1.2 Related works

**Robust Generalization.** One surprising behavior of deep learning is that over-parameterized neural networks can generalize well despite their ability to fit random data (Zhang et al., 2017; Belkin et al., 2019). However, in contrast to the standard (non-robust) generalization, for the robust setting, Rice et al. (2020) empirically investigates robust performance of models based on adversarial training methods, which are designed to improve adversarial robustness (Szegedy et al., 2013; Madry et al., 2017), and shows that *robust overfitting* can be observed on multiple datasets. From the theoretical side, Madry et al. (2017) proposes the notion of robust test error to measure the performance of a model under adversarial attacks, and the required sample complexity is studied in various settings (Schmidt et al., 2018; Bhagoji et al., 2019; Dan et al., 2020; Bhattacharjee et al., 2021). In this paper, we mainly focus on this robust generalization gap and provide a theoretical understanding from the perspective of expressive power.

**Robust interpolation.** Bubeck et al. (2021) proposes a conjecture that over-parameterization is necessary for smooth interpolation. Then Bubeck and Sellke (2021) establishes a law of robustness for isoperimetric data. Specifically, they prove an  $\Omega(\sqrt{Nd/p})$  Lipschitzness lower bound for smooth interpolation, where  $N$ ,  $d$ , and  $p$  denote the sample size, the inputs’ dimension, and the number of parameters, respectively. This result indicates that over-parameterization may be necessary for robust learning. Zhang et al. (2022) studies many data are needed for robust interpolation. This line

of works focuses on the training error and the worst-case robustness (*i.e.* Lipschitz constant), while we measure robustness via the robust generalization error.

**Memorization power of neural networks.** Our work is related to another line of works (e.g., Baum, 1988; Yun et al., 2019; Bubeck et al., 2020; Zhang et al., 2021a; Rajput et al., 2021; Zhang et al., 2021b; Vardi et al., 2021) on the memorization power of neural networks. Among these works, Yun et al. (2019) shows that a neural network with  $\mathcal{O}(N)$  parameters can memorize the data set with zero error, where  $N$  is the size of the data set. Under an additional separable assumption, Vardi et al. (2021) derives an improved upper bound of  $\tilde{\mathcal{O}}(\sqrt{N})$ , which is shown to be optimal. In this work, we show that  $\tilde{\mathcal{O}}(Nd)$  parameters is sufficient for achieving low robust training error. This is in contrast with our exponential lower bound for low robust *test* error.

**Function approximation.** Our work is related to a line of works on function approximation via neural networks (e.g., Cybenko, 1989; Hornik, 1991; Lu et al., 2017; Yarotsky, 2017; Hanin, 2019). Yarotsky (2017) is the most related, which shows that the functions in Sobolev spaces can be uniformly approximated by deep ReLU networks. Also related is the studies of using deep ReLU networks to approximate functions supported on low dimensional manifolds (Chui and Mhaskar, 2018; Shaham et al., 2018; Chen et al., 2019). In particular, Chen et al. (2019) proves that any  $C^n$  function in Hölder spaces can be  $\epsilon$ -approximated by the neural network with size  $\mathcal{O}(\epsilon^{-k/n})$ , where  $k$  is the intrinsic dimension of the manifold embedded in  $\mathbb{R}^d$ . In the robust classification scenario, we can also achieve dimensionality reduction for low-dimensional data.

### 1.3 Notations

Throughout this paper, we use  $\|\cdot\|_p, p \in [1, +\infty]$  to denote the  $\ell_p$  norm in the vector space  $\mathbb{R}^d$ . For  $\mathbf{x} \in \mathbb{R}^d$  and  $A \subset \mathbb{R}^d$ , we can define the distance between  $\mathbf{x}$  and  $A$  as  $d_p(\mathbf{x}, A) = \inf\{\|\mathbf{x} - \mathbf{y}\|_p : \mathbf{y} \in A\}$ . For  $r > 0$ ,  $\mathcal{B}_p(\mathbf{x}, r) = \{\mathbf{y} \in \mathbb{R}^d : \|\mathbf{x} - \mathbf{y}\|_p \leq r\}$  is defined as the  $\ell_p$ -ball with radius  $r$  centered at  $\mathbf{x}$ . For a function class  $\mathcal{F}$ , we use  $d_{VC}(\mathcal{F})$  to denote its VC-dimension. A *multilayer neural network* is a function from input  $\mathbf{x} \in \mathbb{R}^d$  to output  $\mathbf{y} \in \mathbb{R}^m$ , recursively defined as follows:

$$\begin{aligned} \mathbf{h}_1 &= \sigma(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1), \quad \mathbf{W}_1 \in \mathbb{R}^{m_1 \times d}, \mathbf{b}_1 \in \mathbb{R}^{m_1}, \\ \mathbf{h}_\ell &= \sigma(\mathbf{W}_\ell \mathbf{h}_{\ell-1} + \mathbf{b}_\ell), \quad \mathbf{W}_\ell \in \mathbb{R}^{m_\ell \times m_{\ell-1}}, \mathbf{b}_\ell \in \mathbb{R}^{m_\ell}, 2 \leq \ell \leq L-1, \\ \mathbf{y} &= \mathbf{W}_L \mathbf{h}_L + \mathbf{b}_L, \quad \mathbf{W}_L \in \mathbb{R}^{m \times m_L}, \mathbf{b}_L \in \mathbb{R}^m, \end{aligned}$$

where  $\sigma$  is the activation function and  $L$  is the depth of the neural network. In this paper, we mainly focus on ReLU networks *i.e.*  $\sigma(x) = \max\{0, x\}$ . The size of a neural network is defined as its number of weights/parameters *i.e.* the number of its non-zero connections between layers.

## 2 Mild Over-parameterized ReLU Nets Achieve Zero Robust Training Error

With access to only finite amount of data, a common practice for learning a robust classifier is to minimize the *robust training error* (defined below). In this section, we show that neural networks with reasonable size can achieve zero robust training error on a finite training set.

**Definition 2.1** (Robust training error). *Given a data set  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{1 \leq i \leq N}$ ,  $y_i \in \{-1, +1\}$  and an adversarial perturbation radius  $\delta \geq 0$ , the robust training error of a classifier  $f$  is defined as  $\hat{\mathcal{L}}_{\mathcal{D}}^{p, \delta}(f) = \frac{1}{N} \sum_{i=1}^N \mathbb{I}\{\exists \mathbf{x}' \in \mathcal{B}_p(\mathbf{x}_i; \delta), \text{sgn}(f(\mathbf{x}')) \neq y_i\}$ .*

When  $\delta = 0$ , the definition coincides with the standard training error. In this paper, we mainly focus on the case  $p = 2$  and  $p = \infty$ , but our results can be extended to general  $p$  as well.

The following is our main result in this section, which states that for binary classification problems, a neural network with  $\tilde{\mathcal{O}}(Nd)$  weights can perfectly achieve robust classification on a data set of size  $N$ . The detailed proof is deferred to Appendix B.3.

**Theorem 2.2.** *Suppose that  $\mathcal{D} \subset \mathcal{B}_p(0, 1)$  with  $p \in \{2, +\infty\}$  consists of  $N$  data, and the two classes in  $\mathcal{D}$  are  $2\epsilon$ -separated (cf. Definition 1.1), where  $\epsilon \in (0, \frac{1}{2})$  is a constant. Let robustness radius  $\delta < \frac{1}{2}\epsilon$ , then there exists a classifier  $f$  represented by a ReLU network with at most*

$$\mathcal{O}(Nd \log(\delta^{-1}d) + N \cdot \text{polylog}(\delta^{-1}N))$$

parameters, such that  $\hat{\mathcal{L}}_{\mathcal{D}}^{p,\delta}(f) = 0$ .

Theorem 2.2 implies that neural networks is quite efficient for robust classification of finite training data. We also derive a lower bound in the same setting, which is stated below. It is an interesting future direction to study whether this lower bound can be achieved.

**Theorem 2.3.** *Let  $p \in \{2, +\infty\}$  and  $\mathcal{F}_n$  be the set of functions represented by ReLU networks with at most  $n$  parameters. For arbitrary  $2\epsilon$ -separated data set  $\mathcal{D}$  under  $\ell_p$  norm, if there exists a classifier  $f \in \mathcal{F}_n$  such that  $\hat{\mathcal{L}}_{\mathcal{D}}^{p,\delta}(f) = 0$ , then it must hold that  $n = \Omega(\sqrt{Nd})$ .*

The detailed proof of Theorem 2.3 is in Appendix B.4. We leave it as a future direction to study whether this lower bound can be attained. While optimal (non-robust) memorization of  $N$  data points only needs constant width (Vardi et al., 2021), our construction in Theorem 2.2 has width  $\tilde{\mathcal{O}}(Nd)$ . Therefore, if our upper bound is tight, then Theorem 2.2 can probably explain why increasing the network width can benefit robust training (Madry et al., 2017).

### 3 Hardness of Robust Generalization : A Warm-up

In the previous section, we give an upper bound on the size of ReLU networks to robustly classify finite training data. However, it says nothing about *generalization*, or the robust test error, which is arguably a crucial aspect of evaluating the performance of a trained model. As a warm-up, in this section we first consider the same setting as Section 2 where we only assume the data to be well-separated. We show that in this setting, even achieving high standard test accuracy requires exponentially large neural networks in the worst case, which is quite different from empirical observations. This motivates to consider data with additional structures in subsequent sections.

**Definition 3.1** (Robust test error). *Given a probability measure  $P$  on  $\mathbb{R}^d \times \{-1, +1\}$  and a robust radius  $\delta \geq 0$ , the robust test error of a classifier  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  w.r.t  $P$  and  $\delta$  under  $\ell_p$  norm is defined as  $\mathcal{L}_P^{p,\delta}(f) = \mathbb{E}_{(\mathbf{x}, y) \sim P} [\max_{\|\mathbf{x}' - \mathbf{x}\|_p \leq \delta} \mathbb{I}\{\text{sgn}(f(\mathbf{x}')) \neq y\}]$ .*

In contrast with the training set which only consists of finite data points, when studying generalization, we must consider potentially infinite points in the classes that we need to classify. As a result, we consider two disjoint sets  $A, B \in [0, 1]^d$ , where points in  $A$  have label  $+1$  and points in  $B$  have label  $-1$ . We are interested in the following questions:

- Does there exists a robust classifier of  $A$  and  $B$ ?
- If so, can we derive upper and lower bounds on the size of a neural network to robustly classify  $A$  and  $B$ ?

It turns out that, similar to the previous section, the  $\epsilon$ -separated condition (cf. Definition 1.1) ensures the existence of such a classifier. Moreover, it can be realized by a Lipschitz function. This fact has been observed in Yang et al. (2020), and we provide a different version of their result below for completeness.

**Proposition 3.2.** *For  $2\epsilon$ -separated  $A, B \subset [0, 1]^d$  under  $\ell_p$  norm with  $p \in \{2, +\infty\}$ , the classifier  $f^*(\mathbf{x}) := \frac{d_p(\mathbf{x}, B) - d_p(\mathbf{x}, A)}{d_p(\mathbf{x}, A) + d_p(\mathbf{x}, B)}$  is  $\epsilon^{-1}$ -Lipschitz continuous, and satisfies  $\mathcal{L}_P^{p,\epsilon}(f^*) = 0$  for any probability distribution  $P$  on  $A \cup B$ .*

Based on this observation, Yang et al. (2020) concludes that adversarial training is not inherently hard. Rather, they argue that current pessimistic results on robust test error is due to the limits of existing algorithms. However, it remains unclear whether the Lipschitz function constructed in Proposition 3.2 can actually be efficiently approximated by neural networks. The following theorem shows that ReLU networks with exponential size is sufficient for as robust classification.

**Theorem 3.3.** *For any two  $2\epsilon$ -separated  $A, B \subset [0, 1]^d$  under  $\ell_p$  norm with  $p \in \{2, +\infty\}$ , distribution  $P$  on the supporting set  $S = A \cup B$  and robust radius  $c \in (0, 1)$ , there exists a ReLU network  $f$  with at most  $\tilde{\mathcal{O}}(((1 - c)\epsilon)^{-d})$  parameters, such that  $\mathcal{L}_P^{p,c\epsilon}(f) = 0$ .*

The detailed proof is deferred to Appendix C.1. Indeed, it is well known that without additional assumptions, an exponentially large number of parameters is also *necessary* for approximating a Lipschitz function (DeVore et al., 1989; Shen et al., 2022). This result motivates us to consider the

second question listed above. The following result implies that even *without* requiring robustness, neural networks need to be exponentially large to correctly classify  $A$  and  $B$ :

**Theorem 3.4.** *Let  $\mathcal{F}_n$  be the set of functions represented by ReLU networks with at most  $n$  parameters. Suppose that for any  $2\epsilon$ -separated sets  $A, B \subset [0, 1]^d$  under  $\ell_p$  norm with  $p \in \{2, +\infty\}$ , there exists  $f \in \mathcal{F}_n$  that can classify  $A, B$  with zero (standard) test error, then it must hold that  $n = \Omega\left((2\epsilon)^{-\frac{d}{2}} (d \log(1/2\epsilon))^{-\frac{1}{2}}\right)$ .*

Theorem 3.4 implies that mere separability of data sets is insufficient to guarantee that they can be classified by ReLU networks, unless the network size is exponentially large. The detailed proof is in Appendix C.2.

However, one should be careful when interpreting the conclusion of Theorem 3.4, since real-world data sets may possess additional structural properties. Theorem 3.4 does not take these properties into consideration, so it does not rule out the possibility that these additional properties make robust classification possible. Specifically, the joint distribution of data can be decomposed as

$$\mathcal{P}(X, Y) = \underbrace{\mathcal{P}(Y | X)}_{\text{labeling mapping}} \underbrace{\mathcal{P}(X)}_{\text{input}},$$

where  $\mathcal{P}(X, Y)$ ,  $\mathcal{P}(Y | X)$ , and  $\mathcal{P}(X)$  denote the joint, conditional and marginal distributions, respectively. In subsequent sections, we consider two well-known properties of data sets that correspond to the labeling mapping structure (Section 4) and the input structure (Section 5), respectively, and study whether they can bring improvement to neural networks' efficiency for robust classification.

## 4 Robust Generalization for Linear Separable Data

We have seen that for separated data, if no other structural properties are taken into consideration, even standard generalization requires exponentially large neural networks. However, in practice it is often possible to train neural networks that can achieve fairly high standard test accuracy, indicating a gap between the setting of Section 3 and practice.

This motivates us to consider the following question: assuming that there exists a simple classifier that achieves zero standard test error on the data, is it guaranteed that neural networks with reasonable size can also achieve high *robust* test accuracy?

We give a negative answer to this question. Namely, we show that even in the arguably simplest setting where the given data is linear separable and well-separated (cf. Definition 1.1), ReLU networks still need to be exponentially large to achieve high robust test accuracy.

### 4.1 Main results under the linear separable assumption

Clearly, the robust test error (cf. Definition 3.1) depends on the underlying distribution  $P$ . We consider a class of data distributions which have bounded density ratio with the uniform distribution:

**Definition 4.1** (Balanced distribution). *Let  $S \subset \mathbb{R}^n$  such that there exists a uniform probability measure  $m_0$  on  $S$ . A distribution  $P$  on  $S$  is called  $\mu$ -balanced if*

$$\inf \left\{ \frac{P(E)}{m_0(E)} : E \text{ is Lebesgue measurable and } m_0(E) > 0 \right\} \geq \mu.$$

**Remark 4.2.** *Definition 4.1 has also appeared in (Shafahi et al., 2018, Theorem 1), which gives an impossibility result on robust classification, albeit in a completely different setting. Intuitively, it rules out the possibility that data points in certain regions are heavily under-represented.*

The following theorem is the main result of this paper, and the proof sketch is deferred to Section 4.3.

**Theorem 4.3.** *Let  $\epsilon \in (0, 1)$  be a small constant,  $p \in \{2, +\infty\}$  and  $\mathcal{F}_n$  be the set of functions represented by ReLU networks with at most  $n$  parameters. There exists a sequence  $N_d = \Omega\left((2\epsilon)^{-\frac{d-1}{6}}\right)$ ,  $d \geq 1$  and a universal constant  $C_1 > 0$  such that the following holds: for any  $c \in (0, 1)$ , there exists two linear separable sets  $A, B \subset [0, 1]^d$  that are  $2\epsilon$ -separated under  $\ell_p$  norm,*

such that for any  $\mu_0$ -balanced distribution  $P$  on the supporting set  $S = A \cup B$  and robust radius  $c\epsilon$  we have

$$\inf \{ \mathcal{L}_P^{p, c\epsilon}(f) : f \in \mathcal{F}_{N_d} \} \geq C_1 \mu_0.$$

Theorem 4.3 states that the robust test error is lower-bounded by a positive constant  $\alpha = C_1 \mu_0$  unless the ReLU network has size larger than  $\exp(\Omega(d))$ . On the contrary, if we do not require robustness, then the data can be classified by a simple linear function. Moreover, this classifier can be learned with a poly-time efficient algorithm (The detailed proof is in Appendix D.2) :

**Theorem 4.4.** *For any two linear-separable  $A, B \subset [0, 1]^d$ , a distribution  $P$  on the supporting set  $S = A \cup B$ ,  $\delta > 0$  and  $\beta > 0$ , let  $H$  be the family of  $d$ -dimensional hyperplane classifiers. Then, there exists a poly-time efficient algorithm  $\mathcal{A} : 2^S \rightarrow H$ , for  $N = \Omega(d/\beta^2)$  training instances independently randomly sampled from  $P$ , with probability  $1 - \delta$  over samples, we can use the algorithm  $\mathcal{A}$  to learn a classifier  $\hat{f} \in F$  such that*

$$\mathcal{L}_P(\hat{f}) \leq \beta,$$

where  $\mathcal{L}_P(f) := \mathbb{P}_{(\mathbf{x}, y) \sim P} \{y \neq f(\mathbf{x})\}$  denotes the standard test error.

The practical implication of Theorem 4.3 is two-fold. First, by comparing with Theorem 4.4, one can conclude that robust classification may require exponentially more parameters than the non-robust case, which is consistent with the common practice that larger models are used for adversarial robust training. Second, together with our upper bound in Theorem 2.2, Theorem 4.3 implies an exponential separation of neural network size for achieving high robust training and test accuracy.

## 4.2 Exponential lower bound for more general models

In general, our lower bounds hold true for a variety of neural network families and other function classes as well, as long as their VC dimension is at most polynomial in the number of parameters, which is formally stated as Theorem 4.5 that can be derived by the proof of Theorem 4.3 directly.

**Theorem 4.5.** *Let  $\epsilon \in (0, 1)$  be a small constant,  $p \in \{2, +\infty\}$  and  $\mathcal{G}_n$  be the family of parameterized models with at most  $n$  parameters, satisfying the VC-dimension of function family  $\text{VC-dim}(\mathcal{G}_n)$  is at most  $\text{poly}(n)$ . Then, there exists a sequence  $N_d = \exp(\Omega(d))$ ,  $d \geq 1$  and a universal constant  $C'_1 > 0$  such that the following holds: for any  $c \in (0, 1)$ , there exists two linear separable sets  $A, B \subset [0, 1]^d$  that are  $2\epsilon$ -separated under  $\ell_p$  norm, such that for any  $\mu_0$ -balanced distribution  $P$  on the supporting set  $S = A \cup B$  and robust radius  $c\epsilon$  we have*

$$\inf \{ \mathcal{L}_P^{p, c\epsilon}(g) : g \in \mathcal{G}_{N_d} \} \geq C'_1 \mu_0.$$

In other words, the robust generalization error cannot be lower than a constant  $\alpha = C'_1 \mu_0$  unless the model, satisfying the property of their VC dimension polynomially bounded by the number of parameters, has exponential larger size. Indeed, this property is satisfied for e.g. feedforward neural networks with sigmoid (Karpinski and Macintyre, 1995) and piecewise polynomial (Bartlett et al., 2019) activation functions. Therefore, our results reveal that the hardness of robust generalization may stem from the expressive power of generally practical models.

## 4.3 Proof sketch of Theorem 4.3

In this subsection, we present a proof sketch for Theorem 4.3 in the  $\ell_\infty$ -norm case. We only consider  $P$  to be the uniform distribution, extending to  $\mu_0$ -balanced distributions is not difficult, The case of  $\ell_2$ -norm is similar and can be found in the Appendix.

*Proof Sketch.* Let  $K = \lfloor \frac{1}{2\epsilon} \rfloor$ , and  $\phi : \{1, 2, \dots, K\}^{d-1} \rightarrow \{-1, +1\}$  be an arbitrary mapping, we define  $S_\phi = \left\{ \left( \frac{i_1}{K}, \frac{i_2}{K}, \dots, \frac{i_{d-1}}{K}, \frac{1}{2} + \epsilon_0 \cdot \phi(i_1, i_2, \dots, i_{d-1}) \right) : 1 \leq i_1, i_2, \dots, i_{d-1} \leq K \right\}$ , where  $\epsilon_0$  is an arbitrarily small constant. The hyperplane  $x^{(d)} = \frac{1}{2}$  partitions  $S_\phi$  into two subsets, which we denote by  $A_\phi$  and  $B_\phi$ . It is not difficult to check that  $A_\phi$  and  $B_\phi$  satisfies all the required conditions.

Our goal is to show that there exists some choice of  $\phi$  such that robust classification is hard. To begin with, suppose that robust classification with accuracy  $1 - \alpha$  can be achieved with at most  $M$

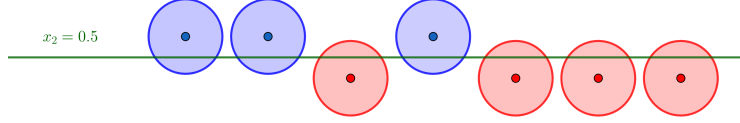


Figure 1: An example of our construction for  $d = 2$ . We choose  $A, B$  as the set of blue points and red points, respectively.

parameters for all  $\phi$ , then these networks can all be embedded into an *enveloping network*  $F_\theta$  of size  $\mathcal{O}(M^3)$ .

Define  $\tilde{S} = \left\{ \left( \frac{i_1}{K}, \frac{i_2}{K}, \dots, \frac{i_{d-1}}{K}, \frac{1}{2} \right) : 1 \leq i_1, i_2, \dots, i_{d-1} \leq K \right\}$ . Robustness implies that for all possible label assignment to  $\tilde{S}$ , at least  $(1 - \alpha)K^{d-1}$  points can be correctly classified by  $F_\theta$ .

If  $\alpha = 0$  i.e. perfect classification is required, then we can see that the  $\text{VC-dim}(F_\theta) \geq K^{d-1}$ , which implies that its size must be exponential, by applying classical VC-dimension upper bounds of neural networks (Bartlett et al., 2019).

When  $\alpha > 0$ , we cannot directly use the bound on VC-dimension. Instead, we use a double-counting argument to lower-bound the *growth function* of some subset of  $\tilde{S}$ .

Let  $V = \frac{1}{2}K^{d-1}$ . Each choice of  $\phi$  corresponds to  $\binom{(1-\alpha)K^{d-1}}{V}$  labelled  $V$ -subset of  $\tilde{S}$  that are correctly classified. There are a total of  $2^{K^{d-1}}$  choices of  $\phi$ , while each labelled  $V$ -subset can be obtained by at most  $2^{K^{d-1}-V}$  different  $\phi$ . As a result, the total number of labelled  $V$ -subset correctly classified by  $F_\theta$  is at least  $2^V \binom{(1-\alpha)K^{d-1}}{V}$ .

On the other hand, the total number of  $V$ -subset of  $\tilde{S}$  is  $\binom{K^{d-1}}{V}$ , thus there must exists a  $V$ -subset  $\mathcal{V}_0 \subset \tilde{S}$ , such that at least

$$\binom{K^{d-1}}{V}^{-1} \cdot 2^V \binom{(1-\alpha)K^{d-1}}{V} \geq \left( \frac{2 \binom{(1-\alpha)K^{d-1} - V}{V}}{K^{d-1} - V} \right)^V \geq C_\alpha^{K^{d-1}} \quad (1)$$

different labellings of  $\mathcal{V}_0$  are correctly classified by  $F_\theta$ , where  $C_\alpha = \sqrt{2(1-2\alpha)} > 1$  for  $\alpha = 0.1$ . Since (1) provides a lower bound for the growth function, together with the upper bound of growth function in terms of VC-dimension, we can deduce that  $\text{VC-dim}(F_\theta) \geq 0.05K^{d-1}$ . Finally, the conclusion follows by applying the VC-dimension bounds in Bartlett et al. (2019).  $\square$

**Remark 4.6.** *The connection between VC dimension and approximation error has been explored in a number of previous works (Yarotsky, 2017; Shen et al., 2022) to provide lower bounds on the network size for approximating a given function class. Here we consider the problem of robust classification which is of more practical interest than function approximation, and our main technical contribution is an exponential lower bound on the VC dimension. Our proof formalizes the folklore that adversarial training is hard since it requires a more complicated decision boundary. We note that similar ideas have been used to show benefits of depth in neural networks (Telgarsky, 2016; Liang and Srikant, 2017) but their techniques are restricted to one-dimensional functions.*

## 5 Robust Generalization for Low-Dimensional-Manifold Data

In this section, we focus on refined structure of data's input distribution  $\mathcal{P}(X)$ . A common belief of real-life data such as images is that the data points lie on a low-dimensional manifold. It promotes a series of methods that are invented to make the dimensionality reduction, including linear dimensionality reduction (e.g., PCA (Pearson, 1901)) and non-linear dimensionality reduction (e.g.,  $t$ -SNE (Hinton and van der Maaten, 2008)). Several works have also empirically verified the belief. Roweis and Saul (2000) and Tenenbaum et al. (2000) have demonstrated that image, speech and other variant form data can be modeled nearly on low-dimensional manifolds. In particular, Wang et al. (2016) studies auto-encoder based dimensionality reduction, and shows that the  $28 \times 28 = 784$  dimensional image from MNIST can be reduced to nearly 10 dimensional representations, which corresponds to the intrinsic dimension of the handwritten digital dataset.



Motivated by these findings, in this section, we assume that data lies on a low-dimensional manifold  $\mathcal{M}$  embedded in  $[0, 1]^d$  i.e.  $\text{supp}(X) \subset \mathcal{M} \subset [0, 1]^d$ . We will show a improved upper bound that is exponential in the intrinsic  $k$  of the manifold  $\mathcal{M}$  instead of the ambient data dimension  $d$  for the size of networks achieving zero robust test error, which implies the efficiency of robust classification under the low-dimensional manifold assumption. Also, we point out that the exponential dependence of  $k$  is not improvable by establishing a matching lower bound.

## 5.1 Preliminaries

Let  $\mathcal{M}$  be a  $k$ -dimensional compact Riemannian manifold embedded in  $\mathbb{R}^d$ , where  $k$  is the intrinsic dimension ( $k \ll d$ ).

**Definition 5.1** (Chart, atlas and smooth manifold). *A chart for  $\mathcal{M}$  is a pair  $(U, \phi)$  such that  $U \subset \mathcal{M}$  is open and  $\phi : U \rightarrow \mathbb{R}^k$ , where  $\phi$  is a homeomorphism; An atlas for  $\mathcal{M}$  is a collection  $\{(U_\alpha, \phi_\alpha)\}_{\alpha \in A}$  of pairwise  $C^n$  compatible charts such that  $\bigcup_{\alpha \in A} U_\alpha = \mathcal{M}$ ; And we call  $\mathcal{M}$  a smooth manifold if and only if  $\mathcal{M}$  has a  $C^\infty$  atlas.*

**Definition 5.2** (Partition of unity). *A  $C^\infty$  partition of unity on a manifold  $\mathcal{M}$  is a collection of non-negative  $C^\infty$  functions  $\rho_\alpha : \mathcal{M} \rightarrow \mathbb{R}_+$  for  $\alpha \in A$  that satisfy (1) the collection of supports,  $\{\text{supp}(\rho_\alpha)\}_{\alpha \in A}$ , is locally finite; and (2)  $\sum_{\alpha \in A} \rho_\alpha = 1$ .*

**Definition 5.3** (Poly-Partitionable). *We call that  $\mathcal{M}$  is poly-partitionable if and only if, for a tangent-space-induced atlas  $\{(U_\alpha, T_\alpha)\}_{\alpha \in A}$  of  $\mathcal{M}$ , there exists a particular partition of unity  $\{\rho_\alpha\}_{\alpha \in A}$  that satisfies  $\rho_\alpha \circ T_\alpha^{-1}$  is a simple piecewise polynomial in  $\mathbb{R}^k$ , where simple piecewise polynomial is defined as the composite mapping between a polynomial and a size-bounded ReLU network.*

The concept, poly-partitionable, defines a class of manifolds that have simple partition of unity, which is a generalization of some structures in the standard Euclidean space  $\mathbb{R}^d$ . For example, an explicit construction for low-dimensional manifold  $[0, 1]^k$  is  $\{\phi_m(x)\}$  in Yarotsky (2017), where the coordinate system is identity mapping.

## 5.2 Main results under the low-dimensional manifold assumption

Before giving our main results, we first extend robust classification to the version of manifold.

**Definition 5.4** (Robust classification on a manifold). *Given a probability measure  $P$  on  $\mathcal{M} \times \{-1, +1\}$  and a robust radius  $\delta$ , the robust test error of a classifier  $f : \mathcal{M} \rightarrow \mathbb{R}$  w.r.t  $P$  and  $\delta$  under  $\ell_p$  norm is defined as  $\mathcal{L}_{\mathcal{M}, P}^{p, \delta}(f) = \mathbb{E}_{(x, y) \sim P} [\max_{x' \in \mathcal{M}, \|x' - x\|_p \leq \delta} \mathbb{I}\{\text{sgn}(f(x')) \neq y\}]$ .*

Now, we present our main result in this section, which establishes an improved upper bound for size that is mainly exponential in the intrinsic dimension  $k$  instead of the ambient data dimension  $d$ .

**Theorem 5.5.** *Let  $\mathcal{M} \subset [0, 1]^d$  be a  $k$ -dimensional compact poly-partitionable Riemannian manifold with the condition number  $\tau > 0$ . For any two  $2\epsilon$ -separated  $A, B \subset \mathcal{M}$  under  $\ell_\infty$  norm, distribution  $P$  on the supporting set  $S = A \cup B$  and robust radius  $c \in (0, 1)$ , there exists a ReLU network  $f$  with at most*

$$\tilde{\mathcal{O}} \left( \left( (1 - c) \epsilon / \sqrt{d} \right)^{-\tilde{k}} \right)$$

*parameters, such that  $\mathcal{L}_{\mathcal{M}, P}^{\infty, c\epsilon}(f) = 0$ , where  $\tilde{k} = \mathcal{O}(k \log d)$  is almost linear with respect to the intrinsic dimension  $k$ , only up to a logarithmic factor.*

*Proof sketch.* The proof idea of Theorem 5.5 has two steps. First, we construct a Lipschitz robust classifier  $f^*$  in Proposition 3.2. Then, we regard  $f^*$  as the target function and use a ReLU network  $f$  to approximate it on the manifold  $\mathcal{M}$ . The following lemma is the key technique that shows we can approximate Lipschitz functions on a manifold by using ReLU networks efficiently.

**Lemma 5.6.** *Let  $\mathcal{M} \subset [0, 1]^d$  be a  $k$ -dimensional compact poly-partitionable Riemannian manifold with the condition number  $\tau > 0$ . For any small  $\delta > 0$  and a  $L$ -lipschitz function  $g : \mathcal{M} \rightarrow \mathbb{R}$ , there exists a function  $\tilde{g}$  implemented by ReLU network with at most  $\tilde{\mathcal{O}} \left( (\sqrt{d}L/\delta)^{-\tilde{k}} \right)$  parameters, such that  $|g - \tilde{g}| < \delta$  for any  $x \in \mathcal{M}$ , where  $\tilde{k}$  is the same as Theorem 5.5.*

By applying the conclusion of Lemma 5.6, we can approximate the  $1/\epsilon$ -Lipschitz function  $f^*$  in Proposition 3.2 via a ReLU network  $f$  with at most  $\tilde{\mathcal{O}}\left(\exp(\tilde{k})\right)$  parameters, such that the uniform approximation error  $\|f - f^*\|_{\ell_\infty(\mathcal{M})}$  at most  $1 - c$ .

Next, we prove the theorem by contradiction. Assume that there exists some perturbed input  $x'$  that is mis-classified and the original input  $x$  is in  $A$ . So we know  $f(x') < 0$  and  $f^*(x) < \epsilon'$ . This implies  $d_\infty(x', A) < d_\infty(x', B) < \frac{1+\epsilon'}{1-\epsilon'}d_\infty(x', A)$ . Combined with  $d_\infty(x', A) + d_\infty(x', B) \geq d_\infty(A, B) \geq 2\epsilon$ , we have  $d_\infty(x', A) > (1 - \epsilon')\epsilon = c\epsilon$ , which is a contradiction.  $\square$

**Remark 5.7.** *Chen et al. (2019) studies network-based approximation on smooth manifolds, and also establishes an  $\mathcal{O}(\delta^{-k})$  bound for the network's size. However, different from their setting where the approximation error  $\delta$  goes to zero, it is reasonable that the separated distance  $\epsilon$  and robust radius  $c$  are constants in our setting. If we simply follow their proofs, we can only obtain the bound  $\mathcal{O}((\delta/C_{\mathcal{M}})^{-k})$  where  $C_{\mathcal{M}}$  also grows exponentially with respect to  $k$ , which further implies that the final result will be roughly  $\exp(\mathcal{O}(k^2))$ . This bound is too loose, especially when  $k \approx \sqrt{d}$ . To this end, we propose a novel approximation framework so as to improve the bound to  $\exp(\mathcal{O}(k))$ , which is presented as Lemma 5.6. And the detailed proof of Lemma 5.6 is deferred to Appendix E.1.*

Although we have shown that robust classification will be more efficient when data lies on a low-dimensional manifold, there is also a curse of dimensionality, i.e., the upper bound for the network's size is almost exponential in the intrinsic dimension  $k$ . The following result shows that the curse of dimension is also inevitable under the low-dimensional manifold assumption.

**Theorem 5.8.** *Let  $\epsilon \in (0, 1)$  be a small constant. There exists a sequence  $\{N_k\}_{k \geq 1}$  that satisfies  $N_k = \Omega\left((2\epsilon\sqrt{d/k})^{-\frac{k}{2}}\right)$ . and a universal constant  $C_1 > 0$  such that the following holds: let  $\mathcal{M} \subset [0, 1]^d$  be a complete and compact  $k$ -dimensional Riemannian manifold with non-negative Ricci curvature, then there exists two  $2\epsilon$ -separated sets  $A, B \subset \mathcal{M}$  under  $\ell_\infty$  norm, such that for any  $\mu_0$ -balanced distribution  $P$  on the supporting set  $S = A \cup B$  and robust radius  $c \in (0, 1)$ , we have  $\inf\{\mathcal{L}_P^{\infty, c\epsilon}(f) : f \in F_{N_k}\} \geq C_1\mu_0$ .*

In other words, the robust test error is lower-bounded by a positive constant  $\alpha = C_1\mu_0$  unless the neural network has size larger than  $\exp(\Omega(k))$ . The detailed proof of Theorem 5.8 is presented in Appendix E.4.

## 6 Conclusion

This paper provides a new theoretical understanding of the gap between the robust training and generalization error. We show that the ReLU networks with reasonable size can robustly classify the finite training data. On the contrary, even with the linear separable and well-separated assumptions, ReLU networks must be exponentially large to achieve low robust generalization error. Finally, we consider the scenario where the data lies on the low dimensional manifold and prove that the ReLU network, with a size exponentially in the intrinsic dimension instead of the inputs' dimension, is sufficient for obtaining low robust generalization error. We believe our work opens up many interesting directions for future work, such as the tight bounds for the robust classification problem, or the reasonable assumptions that permit the polynomial-size ReLU networks to achieve low robust generalization error.

## Acknowledgement

We thank all the anonymous reviewers for their constructive comments. This work is supported by National Science Foundation of China (NSFC62276005), The Major Key Project of PCL (PCL2021A12), Exploratory Research Project of Zhejiang Lab (No. 2022RC0AN02), and Project 2020BD006 supported by PKUBaidu Fund. Binghui Li is partially supported by National Innovation Training Program of China. Jikai Jin is partially supported by the elite undergraduate training program of School of Mathematical Sciences in Peking University.

## References

- Anthony, M., Bartlett, P. L., Bartlett, P. L. et al. (1999). *Neural network learning: Theoretical foundations*, vol. 9. cambridge university press Cambridge.
- Baraniuk, R. G. and Wakin, M. B. (2009). Random projections of smooth manifolds. *Foundations of computational mathematics*, **9** 51–77.
- Bartlett, P. L., Harvey, N., Liaw, C. and Mehrabian, A. (2019). Nearly-tight vc-dimension and pseudodimension bounds for piecewise linear neural networks. *The Journal of Machine Learning Research*, **20** 2285–2301.
- Baum, E. B. (1988). On the capabilities of multilayer perceptrons. *Journal of complexity*, **4** 193–215.
- Belkin, M., Hsu, D., Ma, S. and Mandal, S. (2019). Reconciling modern machine-learning practice and the classical bias–variance trade-off. *Proceedings of the National Academy of Sciences*, **116** 15849–15854.
- Bhagoji, A. N., Cullina, D. and Mittal, P. (2019). Lower bounds on adversarial robustness from optimal transport. *Advances in Neural Information Processing Systems*, **32**.
- Bhattacharjee, R., Jha, S. and Chaudhuri, K. (2021). Sample complexity of robust linear classification on separated data. In *International Conference on Machine Learning*. PMLR.
- Biggio, B., Corona, I., Maiorca, D., Nelson, B., Šrndić, N., Laskov, P., Giacinto, G. and Roli, F. (2013). Evasion attacks against machine learning at test time. In *Joint European conference on machine learning and knowledge discovery in databases*. Springer.
- Bishop, R. L. (1964). A relation between volume, mean curvature and diameter. In *Euclidean Quantum Gravity*. World Scientific, 161–161.
- Boser, B. E., Guyon, I. M. and Vapnik, V. N. (1992). A training algorithm for optimal margin classifiers. In *Proceedings of the fifth annual workshop on Computational learning theory*.
- Bubeck, S., Eldan, R., Lee, Y. T. and Mikulincer, D. (2020). Network size and weights size for memorization with two-layers neural networks. *arXiv preprint arXiv:2006.02855*.
- Bubeck, S., Li, Y. and Nagaraj, D. M. (2021). A law of robustness for two-layers neural networks. In *Conference on Learning Theory*. PMLR.
- Bubeck, S. and Sellke, M. (2021). A universal law of robustness via isoperimetry. *Advances in Neural Information Processing Systems*, **34**.
- Chen, M., Jiang, H., Liao, W. and Zhao, T. (2019). Efficient approximation of deep relu networks for functions on low dimensional manifolds. *Advances in neural information processing systems*, **32**.
- Chui, C. K. and Mhaskar, H. N. (2018). Deep nets for local manifold learning. *Frontiers in Applied Mathematics and Statistics*, **4** 12.
- Cohen, J., Rosenfeld, E. and Kolter, Z. (2019). Certified adversarial robustness via randomized smoothing. In *International Conference on Machine Learning*. PMLR.
- Cybenko, G. (1989). Approximation by superpositions of a sigmoidal function. *Mathematics of control, signals and systems*, **2** 303–314.
- Dan, C., Wei, Y. and Ravikumar, P. (2020). Sharp statistical guarantees for adversarially robust gaussian classification. In *International Conference on Machine Learning*. PMLR.
- Devlin, J., Chang, M.-W., Lee, K. and Toutanova, K. (2018). Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- DeVore, R. A., Howard, R. and Micchelli, C. (1989). Optimal nonlinear approximation. *Manuscripta mathematica*, **63** 469–478.

- Dohmatob, E. (2019). Generalized no free lunch theorem for adversarial robustness. In *International Conference on Machine Learning*. PMLR.
- Gao, R., Cai, T., Li, H., Wang, L., Hsieh, C. J. and Lee, J. D. (2019). Convergence of adversarial training in overparametrized neural networks. *Advances in Neural Information Processing Systems*, **32**.
- Gong, S., Boddeti, V. N. and Jain, A. K. (2019). On the intrinsic dimensionality of image representations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.
- Goodfellow, I. J., Shlens, J. and Szegedy, C. (2014). Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*.
- Hanin, B. (2019). Universal function approximation by deep neural nets with bounded width and relu activations. *Mathematics*, **7** 992.
- Hassani, H. and Javanmard, A. (2022). The curse of overparametrization in adversarial training: Precise analysis of robust generalization for random features regression. *arXiv preprint arXiv:2201.05149*.
- Hinton, G. and van der Maaten, L. (2008). Visualizing data using t-sne journal of machine learning research.
- Hornik, K. (1991). Approximation capabilities of multilayer feedforward networks. *Neural networks*, **4** 251–257.
- Jumper, J., Evans, R., Pritzel, A., Green, T., Figurnov, M., Ronneberger, O., Tunyasuvunakool, K., Bates, R., Žídek, A., Potapenko, A. et al. (2021). Highly accurate protein structure prediction with alphafold. *Nature*, **596** 583–589.
- Karpinski, M. and Macintyre, A. (1995). Polynomial bounds for vc dimension of sigmoidal neural networks. In *Proceedings of the twenty-seventh annual ACM symposium on Theory of computing*.
- Liang, S. and Srikant, R. (2017). Why deep neural networks for function approximation? In *5th International Conference on Learning Representations, ICLR 2017*.
- Lu, Z., Pu, H., Wang, F., Hu, Z. and Wang, L. (2017). The expressive power of neural networks: A view from the width. *Advances in neural information processing systems*, **30**.
- Madry, A., Makelov, A., Schmidt, L., Tsipras, D. and Vladu, A. (2017). Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*.
- Mohri, M., Rostamizadeh, A. and Talwalkar, A. (2018). *Foundations of machine learning*. MIT press.
- Nakkiran, P. (2019). Adversarial robustness may be at odds with simplicity. *arXiv preprint arXiv:1901.00532*.
- Niyogi, P., Smale, S. and Weinberger, S. (2008). Finding the homology of submanifolds with high confidence from random samples. *Discrete & Computational Geometry*, **39** 419–441.
- Pearson, K. (1901). Liii. on lines and planes of closest fit to systems of points in space. *The London, Edinburgh, and Dublin philosophical magazine and journal of science*, **2** 559–572.
- Raghunathan, A., Xie, S. M., Yang, F., Duchi, J. C. and Liang, P. (2019). Adversarial training can hurt generalization. *arXiv preprint arXiv:1906.06032*.
- Rajput, S., Sreenivasan, K., Papailiopoulos, D. and Karbasi, A. (2021). An exponential improvement on the memorization capacity of deep threshold networks. *Advances in Neural Information Processing Systems*, **34**.
- Rice, L., Wong, E. and Kolter, Z. (2020). Overfitting in adversarially robust deep learning. In *International Conference on Machine Learning*. PMLR.

- Roweis, S. T. and Saul, L. K. (2000). Nonlinear dimensionality reduction by locally linear embedding. *science*, **290** 2323–2326.
- Schmidt, L., Santurkar, S., Tsipras, D., Talwar, K. and Madry, A. (2018). Adversarially robust generalization requires more data. *Advances in neural information processing systems*, **31**.
- Shafahi, A., Huang, W. R., Studer, C., Feizi, S. and Goldstein, T. (2018). Are adversarial examples inevitable? *arXiv preprint arXiv:1809.02104*.
- Shafahi, A., Najibi, M., Ghiasi, M. A., Xu, Z., Dickerson, J., Studer, C., Davis, L. S., Taylor, G. and Goldstein, T. (2019). Adversarial training for free! *Advances in Neural Information Processing Systems*, **32**.
- Shaham, U., Cloninger, A. and Coifman, R. R. (2018). Provable approximation properties for deep neural networks. *Applied and Computational Harmonic Analysis*, **44** 537–557.
- Shen, Z., Yang, H. and Zhang, S. (2022). Optimal approximation rate of relu networks in terms of width and depth. *Journal de Mathématiques Pures et Appliquées*, **157** 101–135.
- Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I. and Fergus, R. (2013). Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*.
- Telgarsky, M. (2016). Benefits of depth in neural networks. In *Conference on learning theory*. PMLR.
- Telgarsky, M. (2017). Neural networks and rational functions. In *International Conference on Machine Learning*. PMLR.
- Tenenbaum, J. B., Silva, V. d. and Langford, J. C. (2000). A global geometric framework for nonlinear dimensionality reduction. *science*, **290** 2319–2323.
- Tramèr, F., Kurakin, A., Papernot, N., Goodfellow, I., Boneh, D. and McDaniel, P. (2018). Ensemble adversarial training: Attacks and defenses. In *International Conference on Learning Representations*.
- Tsipras, D., Santurkar, S., Engstrom, L., Turner, A. and Madry, A. (2019). Robustness may be at odds with accuracy. In *International Conference on Learning Representations*. 2019.
- Vardi, G., Yehudai, G. and Shamir, O. (2021). On the optimal memorization power of relu neural networks. *arXiv preprint arXiv:2110.03187*.
- Voulodimos, A., Doulamis, N., Doulamis, A. and Protopapadakis, E. (2018). Deep learning for computer vision: A brief review. *Computational intelligence and neuroscience*, **2018**.
- Wang, Y., Yao, H. and Zhao, S. (2016). Auto-encoder based dimensionality reduction. *Neurocomputing*, **184** 232–242.
- Yang, Y.-Y., Rashtchian, C., Zhang, H., Salakhutdinov, R. R. and Chaudhuri, K. (2020). A closer look at accuracy vs. robustness. *Advances in neural information processing systems*, **33** 8588–8601.
- Yarotsky, D. (2017). Error bounds for approximations with deep relu networks. *Neural Networks*, **94** 103–114.
- Yun, C., Sra, S. and Jadbabaie, A. (2019). Small relu networks are powerful memorizers: a tight analysis of memorization capacity. *Advances in Neural Information Processing Systems*, **32**.
- Zhang, B., Cai, T., Lu, Z., He, D. and Wang, L. (2021a). Towards certifying l-infinity robustness using neural networks with l-inf-dist neurons. In *International Conference on Machine Learning*. PMLR.
- Zhang, C., Bengio, S., Hardt, M., Recht, B. and Vinyals, O. (2017). Understanding deep learning requires rethinking generalization. *iclr. arXiv preprint arXiv:1611.03530*.

- Zhang, H., Wu, Y. and Huang, H. (2022). How many data are needed for robust learning? *arXiv preprint arXiv:2202.11592*.
- Zhang, H., Yu, Y., Jiao, J., Xing, E., El Ghaoui, L. and Jordan, M. (2019). Theoretically principled trade-off between robustness and accuracy. In *International conference on machine learning*. PMLR.
- Zhang, J., Zhang, Y., Hong, M., Sun, R. and Luo, Z.-Q. (2021b). When expressivity meets trainability: Fewer than  $n$  neurons can work. *Advances in Neural Information Processing Systems*, **34** 9167–9180.

## Checklist

1. For all authors...
  - (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]
  - (b) Did you describe the limitations of your work? [Yes]
  - (c) Did you discuss any potential negative societal impacts of your work? [N/A] Explain: our work is seeking to develop a theoretical understanding of the robust generalization in deep learning.
  - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
2. If you are including theoretical results...
  - (a) Did you state the full set of assumptions of all theoretical results? [Yes]
  - (b) Did you include complete proofs of all theoretical results? [Yes]
3. If you ran experiments...
  - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [N/A]
  - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [N/A]
  - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [N/A]
  - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [N/A]
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
  - (a) If your work uses existing assets, did you cite the creators? [N/A]
  - (b) Did you mention the license of the assets? [N/A]
  - (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
  - (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [N/A]
  - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
5. If you used crowdsourcing or conducted research with human subjects...
  - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
  - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
  - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]

## A Preliminaries

In this section, we recall some standard concepts and results in statistical learning theory.

**Definition A.1** (growth function). *Let  $\mathcal{F}$  be a class of functions from  $\mathcal{X} \subset \mathbb{R}^d$  to  $\{-1, +1\}$ . For any integer  $m \geq 0$ , we define the growth function of  $\mathcal{F}$  to be*

$$\Pi_{\mathcal{F}}(m) = \max_{x_i \in \mathcal{X}, 1 \leq i \leq m} |\{(f(x_1), f(x_2), \dots, f(x_m)) : f \in \mathcal{F}\}|.$$

*In particular, if  $|\{(f(x_1), f(x_2), \dots, f(x_m)) : f \in \mathcal{F}\}| = 2^m$ , then  $(x_1, x_2, \dots, x_m)$  is said to be shattered by  $\mathcal{F}$ .*

**Definition A.2** (Vapnik-Chervonenkis dimension). *Let  $\mathcal{F}$  be a class of functions from  $\mathcal{X} \subset \mathbb{R}^d$  to  $\{-1, +1\}$ . The VC-dimension of  $\mathcal{F}$ , denoted by  $\text{VC-dim}(\mathcal{F})$ , is defined as the largest integer  $m \geq 0$  such that  $\Pi_{\mathcal{F}}(m) = 2^m$ . For real-value function class  $\mathcal{H}$ , we define  $\text{VC-dim}(\mathcal{H}) := \text{VC-dim}(\text{sgn}(\mathcal{H}))$ .*

The following result gives a nearly-tight upper bound on the VC-dimension of neural networks.

**Lemma A.3.** (*Bartlett et al., 2019, Theorem 6*) *Consider a ReLU network with  $L$  layers and  $W$  total parameters. Let  $F$  be the set of (real-valued) functions computed by this network. Then we have  $\text{VC-dim}(F) = O(W \log(WL))$ .*

The growth function is connected to the VC-dimension via the following lemma; see e.g. ([Anthony et al., 1999](#), Theorem 7.6).

**Lemma A.4.** *Suppose that  $\text{VC-dim}(\mathcal{F}) = k$ , then  $\Pi_m(\mathcal{F}) \leq \sum_{i=0}^k \binom{m}{i}$ . In particular, we have  $\Pi_m(\mathcal{F}) \leq (em/k)^k$  for all  $m > k + 1$ .*

**Lemma A.5.** (*Mohri et al., 2018, Corollary 3.4*) *Let  $H$  be a family of functions taking values in  $\{-1, +1\}$  with VC-dimension  $k$ . Then, for any  $\delta > 0$ , with probability at least  $1 - \delta$  over  $m$ -samples training dataset  $S$  i.i.d. drawn from the data distribution  $D$ , the following holds for all  $h \in H$ :*

$$\mathcal{L}_D(h) \leq \mathcal{L}_S(h) + \sqrt{\frac{2k \log \frac{em}{k}}{m}} + \sqrt{\frac{\log \frac{1}{\delta}}{2m}},$$

where  $\mathcal{L}_D(h)$  and  $\mathcal{L}_S(h)$  denote the standard test error and training error, respectively.

For deriving upper and lower bounds in the context of  $\ell_2$ -robustness, we also need to introduce the following concepts.

**Definition A.6** ( $\epsilon$ -covering). *Given a set  $\Theta \subset \mathbb{R}^d$ , we say that  $X = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\} \subset \Theta$  is a  $\delta$ -covering of  $\Theta$  if  $\Theta \subset \cup_{i=1}^n \mathcal{B}_2(\mathbf{x}_i, \delta)$ . The covering number  $\mathcal{C}(\Theta, \delta)$  is defined as the minimal size of a  $\delta$ -covering set of  $\Theta$ .*

The following proposition is straightforward from the definition.

**Proposition A.7.** *Let  $\Theta \subset \mathbb{R}^d$  has volume (i.e. Lebesgue measure)  $V$ , then*

$$\mathcal{C}(\Theta, \delta) \geq v_d \cdot \delta^{-d} V,$$

where  $v_d$  is the volume of a  $d$ -dimensional unit ball.

**Definition A.8** ( $\epsilon$ -packing). *Given a set  $\Theta \subset \mathbb{R}^d$ , we say that  $X = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\} \subset \Theta$  is a  $\delta$ -packing of  $\Theta$  if  $\|\mathbf{x}_i - \mathbf{x}_j\|_2 \geq \delta, \forall i \neq j$ . The packing number  $\mathcal{P}(\Theta, \delta)$  is defined as the maximal size of a  $\delta$ -packing set of  $\Theta$ .*

The relationship between the covering and packing number is given by the following result. For completeness, we also provide a simple proof.

**Proposition A.9.** *For any  $\delta \geq 0$ , we have  $\mathcal{P}(\Theta, \delta) \geq \mathcal{C}(\Theta, \delta)$ .*

*Proof.* Consider a maximal packing  $X = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ . Pick any  $\mathbf{x} \in \Theta$ , then there must exist some  $\mathbf{x}_i \in X$  such that  $\|\mathbf{x} - \mathbf{x}_i\|_2 \leq \delta$ ; otherwise,  $X \cup \{\mathbf{x}\}$  is a larger packing set, which contradicts the definition of  $X$ .

Hence it must hold that  $\Theta \subset \cup_{i=1}^n \mathcal{B}_2(\mathbf{x}_i, \delta)$  i.e.  $X$  is a  $\delta$ -covering of  $\Theta$ . The conclusion follows.  $\square$



## B Proofs for Section 2

To prove Theorem 2.2, we first recall some well-known results of neural networks for approximating simple functions.

**Lemma B.1.** *Let  $\varepsilon > 0$ ,  $0 < a < b$  and  $B \geq 1$  be given.*

- (1). (Yarotsky, 2017, Proposition 3) *There exists a function  $\tilde{\times} : [0, B]^2 \rightarrow [0, B^2]$  computed by a ReLU network with  $\mathcal{O}(\log^2(\varepsilon^{-1}B))$  parameters such that*

$$\sup_{x, y \in [0, B]} |\tilde{\times}(x, y) - xy| \leq \varepsilon,$$

and  $\tilde{\times}(x, y) = 0$  if  $xy = 0$ .

- (2). (Telgarsky, 2017, Lemma 3.5) *There exists a function  $R : [a, b] \rightarrow \mathbb{R}_+$  computed by a ReLU network with  $\mathcal{O}(\log^4(a^{-1}b) \log^3(\varepsilon^{-1}b))$  parameters such that  $\sup_{x \in [a, b]} |R(x) - \frac{1}{x}| \leq \varepsilon$ .*

The following lemma establishes uniform approximation of polynomials and is a slight generalization of (Telgarsky, 2017, Lemma 3.4).

**Lemma B.2.** *Let  $\varepsilon \in (0, 1)$ . Suppose that  $P(x) = \sum_{k=1}^s \alpha_k \prod_{i=1}^{r_k} (x_{k,i} - a_{k,i})$  is a polynomial with  $\max_k r_k = r$  and  $\alpha_k, a_{k,i} \in [0, 1], \forall 1 \leq k \leq s, 1 \leq i \leq r_k$ , and  $P(x) \in [-1, +1]$  for  $\forall x \in [0, 1]^d$ . Then there exists a function  $N(x)$  computed by a ReLU network with  $\mathcal{O}(sr \log(\varepsilon^{-1}sr))$  parameters such that  $\sup_{x \in [0, 1]^d} |P(x) - N(x)| \leq \varepsilon$ .*

*Proof.* It suffices to show that each monomial  $P_k(x) = \prod_{i=1}^{r_k} (x_{k,i} - a_{k,i})$  can be  $\varepsilon$ -approximated using  $\mathcal{O}(r \log(\varepsilon^{-1}r))$  parameters. Firstly, we need at most  $r_k \leq r$  parameters to obtain  $x_{k,i} - a_{k,i}, 1 \leq i \leq r_k$  from a linear transformation. We can then apply Lemma B.1 to perform successive multiplication. Note that we still have  $|x_{k,i} - a_{k,i}| \leq 1$ , which can be used to control the cumulative error of  $\tilde{\times}$ .  $\square$

We are now ready to prove Theorem 2.2. For convenience, we restate this theorem below.

**Theorem B.3.** *Suppose that  $\mathcal{D} \subset \mathcal{B}_p(0, 1)$  with  $p \in \{2, +\infty\}$  consists of  $N$  data, and the two classes in  $\mathcal{D}$  are  $2\varepsilon$ -separated (cf. Definition 1.1), where  $\varepsilon \in (0, \frac{1}{2})$  is a constant. Let robustness radius  $\delta < \frac{1}{2}\varepsilon$ , then there exists a classifier  $f$  represented by a ReLU network with at most*

$$\mathcal{O}(Nd \log(\delta^{-1}d) + N \cdot \text{polylog}(\delta^{-1}N))$$

parameters, such that  $\hat{\mathcal{L}}_{\mathcal{D}}^{p, \delta}(f) = 0$ .

*Proof.* (1). The case  $p = 2$ . First, we choose  $C, \varepsilon_1, \varepsilon_2 > 0$  and  $m \in \mathbb{Z}_+$  that satisfy

$$C((\delta^2 + \varepsilon_1)^m + \varepsilon_2) \leq \frac{1}{4} < 4N \leq C((R^2 - \varepsilon_1)^m - \varepsilon_2). \quad (2)$$

These constants will be specified later. Since for  $\forall \mathbf{x}_0 \in [0, 1]^d$ ,  $\mathbf{x} \rightarrow \|\mathbf{x} - \mathbf{x}_0\|^2$  is a polynomial that consists of  $d$  monomials and with degree 2, satisfying the conditions in Lemma B.2, there exists a function  $\phi_1$  computed by a ReLU network with  $\mathcal{O}(d \log(\varepsilon_1^{-1}d))$  parameters such that  $\sup_{\mathbf{x} \in [0, 1]^d} |\phi_1(\mathbf{x}) - \|\mathbf{x} - \mathbf{x}_0\|^2| \leq \varepsilon_1$ . We may further assume that  $\phi([0, 1]^d) \subset [0, 1]$ , or otherwise we can consider  $\sigma(\phi_1(\mathbf{x})) - \sigma(\phi_1(\mathbf{x}) - 1)$  instead.

Applying Lemma B.2 again, we can see that the function  $x \rightarrow x^m$  on  $[0, 1]$  can be approximated with error  $\varepsilon_2$  by a function  $\phi_2$  computed by a ReLU network with  $\mathcal{O}(m \log(\varepsilon_1^{-1}m))$  parameters. Now we can see that  $1 + C \cdot \phi_2 \circ \phi_1$  is computable by a ReLU network and takes value in  $[1, \frac{5}{4}]$  when  $\mathbf{x} \in \mathcal{B}(\mathbf{x}_0, \delta)$  and in  $(4N + 1, C + 1)$  when  $\mathbf{x} \notin \mathcal{B}(\mathbf{x}_0, R)$  (since  $R \leq 1$ ).

The final step is to choose  $\phi_3$  computed by a ReLU network with  $\mathcal{O}(\log^4 C \log^3(NC))$  parameters such that it approximates  $\frac{1}{x}$  on  $[1, C + 1]$  with error  $< \frac{1}{4N}$ . Hence  $\phi_3 \circ (1 + C \cdot \phi_2 \circ \phi_1)$  is larger

than  $\frac{3}{4}$  inside  $\mathcal{B}(\mathbf{x}_0, \delta)$  and smaller than  $\frac{1}{2N}$  outside  $\mathcal{B}(\mathbf{x}_0, R)$ . This construction uses a total of  $\mathcal{O}(W)$  parameters, where

$$W = d \log(\varepsilon_1^{-1}d) + m \log(\varepsilon_2^{-1}m) + \log^4 C \log^3(NC). \quad (3)$$

Finally, we choose

$$\varepsilon_1 = \frac{R\delta(R-\delta)}{R+\delta}, \quad m = \max\left\{1, \log \frac{32N\delta}{R}\right\}, \quad \varepsilon_2 = \frac{1}{33N} \left(\frac{R(R^2+\delta^2)}{R+\delta}\right)^m,$$

and  $C = \frac{4N}{(R^2-\varepsilon_1)^{m-\varepsilon_2}} = \mathcal{O}(N\delta^{-2m})$ , which satisfies (2). Plugging all expressions into (3), we can see that

$$W = \mathcal{O}\left(d\left(\log d + \log \delta^{-1} + \log(R-\delta)^{-1}\right) + \log^7(\delta^{-1}N)\right).$$

We denote this construction by  $\psi(\mathbf{x}; \mathbf{x}_0, \boldsymbol{\theta})$ , where  $\boldsymbol{\theta}$  consists of all parameters. The arguments above show that there exists  $\boldsymbol{\theta} = \boldsymbol{\theta}(\mathbf{x}_0)$  such that  $\psi(\mathbf{x}; \mathbf{x}_0, \boldsymbol{\theta}) > \frac{3}{4}$  when  $\mathbf{x} \in \mathcal{B}(\mathbf{x}_0, \delta)$  and  $\psi(\mathbf{x}; \mathbf{x}_0, \boldsymbol{\theta}) < \frac{1}{2N}$  when  $\mathbf{x} \notin \mathcal{B}(\mathbf{x}_0, R)$ . Consider the function  $\Psi(\mathbf{x}; \boldsymbol{\theta}_{1:N}) = 4 \sum_{i=1}^N \psi(\mathbf{x}; \mathbf{x}_i, \boldsymbol{\theta}_i) - \frac{5}{2}$ . The total number of parameters in  $\Psi$  is  $\tilde{\mathcal{O}}(Nd)$ . Moreover, if we choose  $\boldsymbol{\theta}_i = \boldsymbol{\theta}(\mathbf{x}_i)$  when  $y_i = 1$  and  $\boldsymbol{\theta}_i = 0$  when  $y_i = -1$ , then  $\Psi$  satisfies the condition in Theorem B.3.

(2). *The case  $p = \infty$ .* To obtain the same result under the  $\ell_\infty$  norm, it suffices to construct a neural network with size  $\mathcal{O}(d)$  parameters to represent the function  $x \rightarrow \|x - x_0\|_\infty$ ; the remaining steps are exactly the same with the  $\ell_2$  case.

Let  $x^{(i)}$  denote the  $i$ -th coordinate of  $\mathbf{x}$ , then  $\|\mathbf{x} - \mathbf{x}_0\|_\infty = \max_{1 \leq i \leq d} |x^{(i)} - x_0^{(i)}|$ . Since

$$|a| = \frac{1}{2} (\max\{a, 0\} + \max\{-a, 0\}),$$

we can see that  $x^{(i)} \rightarrow |x^{(i)} - x_0^{(i)}|$  can be represented by a constant-size ReLU network. Moreover, the function  $\max\{a, b\} = \frac{1}{2}(|a+b| + |a-b|)$ , so that the function  $(a_1, a_2, \dots, a_d) \rightarrow \max_{1 \leq i \leq d} a_i$  can be represented with  $\mathcal{O}(d)$  parameters. To summarize,  $\mathbf{x} \rightarrow \|\mathbf{x} - \mathbf{x}_0\|_\infty$  can be represented using a ReLU network of size  $\mathcal{O}(d)$ , as desired.  $\square$

In the following, we prove Theorem 2.3.

**Theorem B.4** (Restatement of Theorem 2.3). *Let  $p \in \{2, +\infty\}$  and  $\mathcal{F}_n$  be the set of functions represented by some ReLU network with at most  $n$  parameters. If for any  $2\varepsilon$ -separated data set  $\mathcal{D}$  under  $\ell_p$  norm, there exists a classifier  $f \in \mathcal{F}_n$  such that  $\hat{\mathcal{L}}_{\mathcal{D}}^{p,\delta}(f) = 0$ , then it must hold that  $n = \Omega(\sqrt{Nd})$ .*

*Proof.* It follows from the assumption that given any data points  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$  which are pairwise  $2\varepsilon$ -separated, there exists  $f \in \mathcal{F}_n$  being able to achieve zero training error for any binary label. It directly follows from (Gao et al., 2019, Theorem 6.1) that

$$\text{VC-dim}(\mathcal{F}_n) = \Omega(Nd).$$

On the other hand, suppose that  $L \neq n$  is the depth of the neural network, then we have

$$\text{VC-dim}(\mathcal{F}_n) = \mathcal{O}(nL \log(nL)) = \mathcal{O}(n^2).$$

As a result, it follows that  $n = \tilde{\Omega}(\sqrt{Nd})$ , as desired.  $\square$

## C Proofs for Section 3

### C.1 Proof of Theorem 3.3

The proof idea of Theorem 3.3 has two key steps. First, we construct a Lipschitz classifier  $f^*$  based on distance function between a point and a close set that can  $\varepsilon$ -robustly classify  $A, B$ . Then we regard  $f^*$  as the target function and use a ReLU network to approximate it to derive the  $c\varepsilon$ -robust classifier. Before proving the theorem, we first introduce the two following useful conclusions, which also corresponding to the two steps of proof.

**Proposition C.1.** For the separable  $A, B \subset [0, 1]^d$ , we define  $f^*(\mathbf{x}) := \frac{d_\infty(\mathbf{x}, B) - d_\infty(\mathbf{x}, A)}{d_\infty(\mathbf{x}, A) + d_\infty(\mathbf{x}, B)}$ , which has the following properties:

1.  $f^*(\mathbf{x})$  can classify  $A, B$  correctly i.e.  $f^*(\mathbf{x}) = \begin{cases} 1, & x \in A \\ -1, & x \in B \end{cases}$ .
2.  $f^*(\mathbf{x})$  is a  $\epsilon$ -robust classifier i.e. for any perturbed input  $\mathbf{x}'$  that satisfies  $\|\mathbf{x}' - \mathbf{x}\|_\infty \leq \epsilon$  can also be classified correctly.
3.  $f^*(\mathbf{x})$  is  $\frac{1}{\epsilon}$ -Lipschitz w.r.t.  $\ell_\infty$  norm.

We can check these properties by the continuity and 1-Lipschitz property of distance function  $d_\infty(p, S)$ .

**Lemma C.2.** For any  $L$ -lipschitz function  $f$  in  $[0, 1]^d$ , there exists a function  $\tilde{f}$  implemented by ReLU network with at most  $c_1(c_2\epsilon/L)^{-d}(d^2 + d \log d + d \log(1/\epsilon))$  parameters that satisfies  $|f(\mathbf{x}) - \tilde{f}(\mathbf{x})| \leq \epsilon$  for any  $\mathbf{x} \in [0, 1]^d$ , where  $c_1$  and  $c_2$  are constants.

This lemma provides a useful approximation tool for us, which is an improved version of Theorem 1 in Yarotsky (2017). Compared with Theorem 1 in Yarotsky (2017), we use Lipschitz property of function instead of high-order differentiability and focus on not the bound order when  $\epsilon$  goes to zero but also more accurate bound order depending on  $\epsilon, L$  and  $d$ . By a refined analysis of total approximation error, we can derive this lemma.

*Proof of Theorem 3.3.* By Lemma C.2, we can approximate  $f^*$  in Lemma C.1 satisfying uniform error at most  $1-c$  via a ReLU network  $f$  with at most  $c_1(c_2(1-c)\epsilon)^{-d}(d^2 + d \log d + d \log(1/(1-c)))$  parameters. Then, we prove the theorem by contradiction. Assume that there exists some perturbed input  $\mathbf{x}'$  that is mis-classified and the original input  $\mathbf{x}$  is in  $A$ . So we know  $f(\mathbf{x}') < 0$  and  $f^*(\mathbf{x}) < \epsilon'$ . This implies  $d_\infty(\mathbf{x}', A) < d_\infty(\mathbf{x}', B) < \frac{1+\epsilon'}{1-\epsilon'}d_\infty(\mathbf{x}', A)$  Combined with  $d_\infty(\mathbf{x}', A) + d_\infty(\mathbf{x}', B) \geq d_\infty(A, B) \geq 2\epsilon$ , we have  $d_\infty(\mathbf{x}', A) > (1 - \epsilon')\epsilon = c\epsilon$ , which is the contradiction.  $\square$

## C.2 Proof of Theorem 3.4

The main idea of proof is to estimate the lower bound of the family's VC-dimension via the definition of  $c\epsilon$ -robust family.

*Proof of Theorem 3.4.* The key idea is to find some discrete points that can be shattered by the function family  $\mathcal{F}_n$ .

(1). *The  $p = \infty$  case.* We use  $K$  to denote  $\lfloor \frac{1}{2\epsilon} \rfloor + 1$ , and we can divide  $[0, 1]^d$  into  $(K - 1)^d$  non-overlapping sub-cubes. Let  $S$  be the set of all the vertices of sub-cubes, which has  $K^d$  elements and can be represented by

$$S = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{K^d}\} = \{(2\epsilon i_1, 2\epsilon i_2, \dots, 2\epsilon i_d) | 0 \leq i_1, i_2, \dots, i_d < K\}.$$

For any partition  $I, J$  of  $[K^d]$  ( $I \cap J = \Phi, I \cup J = [K^d]$ ), let  $A = \{\mathbf{x}_i | i \in I\}$  and  $B = \{\mathbf{x}_j | j \in J\}$  be the positive and negative classes. Then we have  $d_\infty(A, B) \geq 2\epsilon$ . By the definition of  $c\epsilon$ -robust classifier family, there exists a classifier  $f \in \mathcal{F}$  classify  $A, B$  correctly. Thus, the family  $\mathcal{F}$  shatter the subset  $S \subset [0, 1]^d$ . By using the conclusion of Lemma A.3, we have  $K^d \leq \text{VC-dim}(\mathcal{F}) = \mathcal{O}(WL \log(N)) = \mathcal{O}(W^2 \log(W))$  where  $L$  is the depth of networks and  $W$  is the total number of parameters.

(2). *The  $p = 2$  case.* Similar to the case of  $p = \infty$ , we need to construct a set  $S \subset \mathcal{B}_2(0, 1)$  such that the  $\ell_2$ -distance between any two points in  $S$  is as least  $2\epsilon$ .

Specifically, we choose  $S$  to be a  $2\epsilon$ -packing of  $\mathcal{B}_2(0, 1)$  with maximal size. Then we have that  $|S| \geq \mathcal{P}(\mathcal{B}_2(0, 1), 2\epsilon) \geq \mathcal{C}(\mathcal{B}_2(0, 1), 2\epsilon) \geq (2\epsilon)^{-d}$ , by Propositions A.7 and A.9. Similar to the  $p = \infty$  case, robustness implies that  $S$  can be shattered by  $\mathcal{F}_n$ , so that  $K^d = \mathcal{O}(W^2 \log W)$  and the conclusion follows.  $\square$

## D Proofs for Section 4

In this section, we present the proof of Theorem 4.3 and 4.4.

**Theorem D.1** (Restatement of Theorem 4.3). *Let  $\epsilon \in (0, 1)$  be a small constant,  $p \in \{2, \infty\}$  and  $\mathcal{F}_n$  be the set of functions represented by ReLU networks with at most  $n$  parameters. There exists a sequence  $N_d = \Omega\left((2\epsilon)^{-\frac{d-1}{6}}\right)$ ,  $d \geq 1$  and a universal constant  $C > 0$  such that the following holds: for any  $c \in (0, 1)$ , there exists two linear separable sets  $A, B \subset [0, 1]^d$  that are  $2\epsilon$ -separated under  $\ell_p$ -norm, such that for any  $\mu_0$ -balanced distribution  $P$  on the supporting set  $S = A \cup B$  and robust radius  $c\epsilon$  we have*

$$\inf \{ \mathcal{L}_P^{p, c\epsilon}(f) : f \in \mathcal{F}_{N_d} \} \geq C\mu_0.$$

*Proof.* (1). *The  $p = \infty$  case.* Define

$$S_\phi = \left\{ \left( \frac{i_1}{K}, \frac{i_2}{K}, \dots, \frac{i_{d-1}}{K}, \frac{1}{2} + c\epsilon \cdot \phi(i_1, i_2, \dots, i_{d-1}) \right) : 1 \leq i_1, i_2, \dots, i_{d-1} \leq K \right\},$$

and

$$\tilde{S} = \left\{ \left( \frac{i_1}{K}, \frac{i_2}{K}, \dots, \frac{i_{d-1}}{K}, \frac{1}{2} \right) : 1 \leq i_1, i_2, \dots, i_{d-1} \leq K \right\},$$

where  $K = \lfloor \frac{1}{2\epsilon} \rfloor$ , and  $\phi : \{1, 2, \dots, K\}^{d-1} \rightarrow \{-1, +1\}$  be an arbitrary mapping. For a vector  $\mathbf{x} \in \mathbb{R}^d$ , we use  $x^{(i)}$  to denote its  $i$ -th component. Let  $A_\phi = S_\phi \cap \{\mathbf{x} \in \mathbb{R}^d : x^{(d)} > \frac{1}{2}\}$ ,  $B_\phi = S_\phi - A_\phi$  and  $\mu$  be the uniform distribution on  $S$ . It's easy to see that  $A$  and  $B$  are linear separable by the hyperplane  $x^{(d)} = \frac{1}{2}$ . Moreover, we clearly have  $d(A, B) \geq 2\epsilon$ . We will show that there exists some choice of  $\phi$  such that robust classification of  $A_\phi$  and  $B_\phi$  with  $(c\epsilon, 1 - \alpha)$ -accuracy requires at least  $\Omega(K^{(d-1)/6})$  parameters.

Assume that for any choices of  $\phi$ , the induced sets  $A_\phi$  and  $B_\phi$  can always be robustly classified with  $(c\epsilon, 1 - \alpha)$ -accuracy by a ReLU network with at most  $M$  parameters. Then, we can construct an *enveloping network*  $F_\theta$  with  $M - 1$  hidden layers,  $M$  neurons per layer and at most  $M^3$  parameters such that any network with size  $\leq M$  can be embedded into this envelope network. As a result,  $F_\theta$  is capable of  $(c\epsilon, 1 - \alpha)$ -robustly classify any sets  $A_\phi, B_\phi$  induced by arbitrary choices of  $\phi$ . We use  $R_\phi$  to denote the subset of  $S_\phi = A_\phi \cup B_\phi$  satisfying  $|R_\phi| = (1 - \alpha)|S_\phi| = (1 - \alpha)K^{d-1}$  such that  $R_\phi$  can be  $c\epsilon$ -robustly classified.

Consider the projection operator  $\mathcal{P}$  onto the hyperplane  $x^{(d)} = \frac{1}{2}$ . For any set  $C \in \mathbb{R}^d$ , we use  $\tilde{C}$  to denote  $\mathcal{P}(C)$ . Then  $c\epsilon$ -robustness implies that the labelled data set

$$R_\phi^+ = \left\{ (\mathbf{x}, y) : \mathbf{x} \in \tilde{R}_\phi, y = \phi(Kx^{(1)}, \dots, Kx^{(d-1)}) \right\}$$

can be correctly classified by  $F_\theta$ , with appropriate choices of parameters.

Let  $V = \frac{1}{2}K^{d-1}$  and  $\hat{\mathcal{R}}_\phi$  be the collection of all labelled  $V$ -subset (i.e. subset of size  $V$ ) of  $R_\phi^+$ . For each  $V$ -subset  $R$  of  $\tilde{S}$ , we use  $\mathcal{G}(R)$  to denote the set of all labelings of  $R$ , so that  $|\mathcal{G}(R)| = 2^V$ .

Note that for each labelled  $V$ -subset  $T$ , there exists at most  $2^{K^{d-1}-V}$  different choices of  $\phi$  such that  $T \subset R_\phi^+$  (or, equivalently,  $T \in \hat{\mathcal{R}}_\phi$ ): this is because the value of  $\phi$  on data points in  $T$  has been specified by their labels, and there are two choices for each of the remaining  $K^{d-1} - V$  points in  $\{1, 2, \dots, K\}^{d-1}$ . As a result, we have

$$\left| \bigcup_\phi \hat{\mathcal{R}}_\phi \right| \geq 2^{-(K^{d-1}-V)} \sum_\phi |\hat{\mathcal{R}}_\phi| = 2^V \binom{(1-\alpha)K^{d-1}}{V}.$$

On the other hand, the total number of  $V$ -subset of  $\tilde{S}$  is  $\binom{K^{d-1}}{V}$ , thus there must exists a  $V$ -subset  $\mathcal{V}_0 \subset \tilde{S}$ , such that at least

$$\binom{K^{d-1}}{V}^{-1} \cdot 2^V \binom{(1-\alpha)K^{d-1}}{V} \geq \left( \frac{2 \binom{(1-\alpha)k^{d-1} - V}{V}}{K^{d-1} - V} \right)^V \quad (4)$$

different labelings of  $\mathcal{V}_0$  are included in  $\cup_\phi \hat{\mathcal{R}}_\phi$ . Since  $F_\theta$  can correctly classify all elements (which are  $V$ -subsets) in  $\cup_\phi \hat{\mathcal{R}}_\phi$ , it can in particular classify the set  $\mathcal{V}_0$  with at least  $\left(\frac{2((1-\alpha)K^{d-1}-V)}{K^{d-1}-V}\right)^V$  different assignments of labels. Let  $d_{VC}$  be the VC-dimension of  $F_\theta$ , then by Lemma A.4, either  $d_{VC} \geq V = \frac{1}{2}K^{d-1}$ , or

$$(2(1-2\alpha))^V \leq \left(\frac{2((1-\alpha)K^{d-1}-V)}{K^{d-1}-V}\right)^V \leq \Pi_{F_\theta}(V) \leq \left(\frac{eV}{d_{VC}}\right)^{d_{VC}},$$

where  $\Pi$  is the growth function. The RHS is increasing in  $d_{VC}$  as long as  $d_{VC} \leq V$ . When  $\alpha \leq \frac{1}{10}$ , we have  $2(1-2\alpha) > (10e)^{1/10}$ , so that  $d_{VC} \geq \frac{1}{10}V = \frac{1}{20}K^{d-1}$ . Finally, since  $F_\theta$  has at most  $M^3$  parameters, classical bounds on VC-dimension (Bartlett et al., 2019) imply that  $M = \Omega(K^{(d-1)/6})$ , as desired.

(2). *The  $p = 2$  case.* Let  $P$  be an  $2\epsilon$ -packing of the unit ball  $\mathcal{B}_{d-1}$  in  $\mathbb{R}^{d-1}$ . Since the packing number  $\mathcal{P}(\mathcal{B}_{d-1}, \|\cdot\|, 2\epsilon) \geq \mathcal{C}(\mathcal{B}_{d-1}, \|\cdot\|_2, 2\epsilon) \geq (2\epsilon)^{-(d-1)}$  by Propositions A.7 and A.9, where  $\mathcal{C}(\Theta, \|\cdot\|, \epsilon)$  is the  $\epsilon$ -covering number of a set  $\Theta$ . For any  $\lambda \in (0, 1)$ , we can consider the construction

$$S_\phi = \left\{ \left( \mathbf{x}, \frac{1}{2} + \epsilon_0 \cdot \phi(\mathbf{x}) \right) : \mathbf{x} \in P \right\},$$

where  $\phi : P \rightarrow \{-1, +1\}$  is an arbitrary mapping. It's easy to see that all points in  $S_\phi$  with first  $d-1$  components satisfying  $\|\mathbf{x}\|_2 \leq \sqrt{1-\epsilon_0^2}$  are in the unit ball  $\mathcal{B}_d$ , so that by choosing  $\epsilon_0$  sufficiently small, we can guarantee that  $|S_\phi \cap \mathcal{B}_d| \geq \frac{1}{2}(2\epsilon)^{-(d-1)}$ . For convenience we just replace  $S_\phi$  with  $S_\phi \cap \mathcal{B}_d$  from now on.

Let  $A_\phi = S_\phi \cap \{\mathbf{x} \in \mathbb{R}^d : x^{(d)} > \frac{1}{2}\}$ ,  $B_\phi = S_\phi - A_\phi$ . It's easy to see that for arbitrary  $\phi$ , the construction is linear-separable and satisfies  $2\epsilon$ -separability. The remaining steps are just identical to the  $\ell_\infty$  case.  $\square$

**Theorem D.2** (Restatement of Theorem 4.4). *For any two linear-separable  $A, B \subset [0, 1]^d$ , a distribution  $P$  on the supporting set  $S = A \cup B$ ,  $\delta > 0$  and  $\beta > 0$ , let  $H$  be the family of  $d$ -dimensional hyperplane classifiers. Then, there exists a poly-time efficient algorithm  $\mathcal{A} : 2^S \rightarrow H$ , for  $N = \Omega(d/\beta^2)$  training instances independently randomly sampled from  $P$ , with probability  $1 - \delta$  over samples, we can use the algorithm  $\mathcal{A}$  to learn a classifier  $\hat{f} \in F$  such that*

$$\mathcal{L}_P(\hat{f}) \leq \beta,$$

where  $\mathcal{L}_P(f) := \mathbb{P}_{(x,y) \sim P}\{y \neq f(x)\}$  denotes the standard test error.

*Proof.* We i.i.d. sample  $N$  instances from the data distribution  $P$ , and use  $T$  to denote the training dataset. By Lemma A.5, with probability at least  $1 - \delta$ , we have

$$\mathcal{L}_P(h) \leq \mathcal{L}_T(h) + \mathcal{O}\left(\sqrt{\frac{d}{N}}\right), \forall h \in H$$

By conclusions of Boser et al. (1992) and results of convex optimization, we have a poly-time algorithm  $\mathcal{A} : 2^S \rightarrow H$  such that  $\mathcal{L}_T(\mathcal{A}(T)) \leq \frac{\beta}{2}$ , and we use  $\hat{f}$  to denote  $\mathcal{A}(T)$ . Finally, when  $N = \Omega(d/\beta^2)$  is sufficient large, we have  $\mathcal{L}_P(\hat{f}) \leq \frac{\beta}{2} + \frac{\beta}{2} = \beta$ .  $\square$

## E Proofs for Section 5

### E.1 Proof of Lemma 5.6

**Theorem E.1** (Restatement of Lemma 5.6). *Let  $\mathcal{M} \subset [0, 1]^d$  be a  $k$ -dimensional compact poly-partitionable Riemannian manifold with the condition number  $\tau > 0$ . For any small  $\delta > 0$  and a  $L$ -lipschitz function  $f : \mathcal{M} \rightarrow \mathbb{R}$ , there exists a function  $\hat{f}$  implemented by ReLU network with at most*

$$\tilde{\mathcal{O}}\left((\sqrt{d}L/\delta)^{-\tilde{k}}\right)$$

parameters, such that  $|f - \hat{f}| < \delta$  for any  $x \in \mathcal{M}$ , where  $\tilde{k}$  is the same as Theorem 5.5.

*Proof.* The full proof has six steps, and we finally construct a ReLU network as the following form

$$\hat{f} = \sum_{i=1}^N (\hat{f}_i \circ \phi_i) \hat{\times} (\hat{\rho}_i \circ T_i) \hat{\times} (\hat{I}_\theta \circ \hat{d}_i^2),$$

where all  $\hat{f}_i, \phi_i, \hat{\rho}_i, T_i, \hat{I}_\Delta, \hat{d}_i^2$  and  $\hat{\times}$  are implemented by sub ReLU networks, and these notation's detail will be introduced step by step.

As the above form shows, three sub-network groups will be combined by multiplication approximator  $\hat{\times}$ , where each group corresponds to a factor in the partition-of-unity-based decomposition of  $f$  (i.e.  $f = \sum_{i=1}^N f \times \rho_i \times I\{x \in U_i\}$ , and where  $\{\rho_i\}_{i \in [N]}$  satisfying  $\sum_{i \in [N]} \rho_i = 1$  is a partition of unity on an atlas  $\{U_i\}_{i \in [N]}$ ).

### Step 1: Construct poly-partition of unity on $\mathcal{M}$

Consider a open cover  $\{B_r(x)\}_{x \in \mathcal{M}}$  on  $\mathcal{M}$ , where we use  $B_r(c)$  to denote the Euclidean neighborhood with center  $c$  and radius  $r$ . Due to the compactness of manifold  $\mathcal{M}$ , we know there exists a finite open cover  $\{B_r(x_i)\}_{i \in I}$  indexed by a finite sub-index set  $I$ , which satisfies  $\mathcal{M} \subset \bigcup_{i \in I} B_r(x_i)$ .

Then, we estimate the cardinal number of index set. By the conclusions of [Niyogi et al. \(2008\)](#), when we select radius  $r$  satisfying  $r < \tau/2$ , we have the following lemma, which gives an lower bound of  $k$ -dimensional volume of the local neighborhood of  $\mathcal{M}$ .

**Lemma E.2.** ([Niyogi et al., 2008, Lemma 5.3](#)) *Let  $c \in \mathcal{M}$ . Now consider  $U = \mathcal{M} \cap B_r(c)$ . Then  $\text{vol}(U) \geq (\cos(\theta))^k \text{vol}(B_r^k(c))$  where  $B_r^k(c)$  is the  $k$ -dimensional ball in  $T_c$  centered at  $c, \theta = \arcsin(r/2\tau)$ . All volumes are  $k$ -dimensional volumes where  $k$  is the dimension of  $\mathcal{M}$ .*

Recall the relation between the covering number  $\mathcal{N}(\mathcal{M}, d_2, r)$  and the packing number  $\mathcal{P}(\mathcal{M}, d_2, r)$ , and then we have

$$\begin{aligned} \mathcal{N}(\mathcal{M}, d_2, r) &\leq \mathcal{P}(\mathcal{M}, d_2, r/2) \\ &\leq \frac{\text{vol}(\mathcal{M})}{(\cos(\theta))^k \text{vol}(B_{r/2}^k(c))} \\ &\leq c_N \frac{\text{vol}(\mathcal{M})}{r^k}, \end{aligned}$$

where  $c_N$  is a constant that only exponentially depends on  $k \log k$ .

By the poly-partitionable and smoothness properties of Riemannian manifold  $\mathcal{M}$ , there exists a collection  $\{U_i, T_i, \rho_i\}_{i \in \mathcal{N}(\mathcal{M}, d_2, r)}$  such that  $\{U_i, T_i\}$  compose a tangent-space-induced atlas and  $\{\rho_i\}$  also compose a poly-partition of unity on  $\mathcal{M}$ . So we can decompose  $f$  as  $f = \sum_{i=1}^N f \rho_i$ , where we use notation  $N$  to denote  $\mathcal{N}(\mathcal{M}, d_2, r)$  so as to simplify the written process.

### Step 2: Local almost isotropic transformation via random projection

To achieve dimensional reduction of  $f$  in the local neighborhood  $U_i$ , we will use the following random projection technique proposed by [Baraniuk and Wakin \(2009\)](#).

**Lemma E.3.** ([Baraniuk and Wakin, 2009, Theorem 3.1](#))

*Let  $\mathcal{M}$  be a compact  $k$ -dimensional sub-manifold of  $\mathbb{R}^d$  having condition number  $1/\tau$ . Fix  $0 < \delta < 1$  and  $0 < \eta < 1$ . Let  $A$  be a random orthoprojector from  $\mathbb{R}^d$  to  $\mathbb{R}^{\tilde{k}}$  with*

$$\tilde{k} = O\left(\frac{k \log(d \text{vol}(\mathcal{M}) \tau^{-1} \delta^{-1}) \log(1/\eta)}{\delta^2}\right).$$

*If  $\tilde{k} \leq d$ , then with probability at least  $1 - \eta$  the following statement holds: For every distinct pair of points  $x, y \in \mathcal{M}$ ,*

$$(1 - \delta) \sqrt{\frac{\tilde{k}}{d}} \leq \frac{\|Ax - Ay\|_2}{\|x - y\|_2} \leq (1 + \delta) \sqrt{\frac{\tilde{k}}{d}}.$$

Since we can select  $\eta$  is very close to 1 in order to the probability  $1 - \eta > 0$ , there exists a orthoprojector  $A_i$  for sub-manifold  $U_i$  by applying Lemma E.3. And we use  $V_r$  to denote the uniform upper bound of  $\text{vol}(U_i)$ , which makes the uniform dimension  $\tilde{k} = O(k \log d)$  for each  $i \in [N]$ . Let local almost isotropic transformation  $\phi_i(x) = \frac{1}{2}A_i(x - c_i) + \frac{1}{2}\mathcal{K}$ , where we use  $\mathcal{K}$  to denote the vector  $(1, 1, \dots, 1) \in \mathbb{R}^{\tilde{k}}$ , and then we know  $\phi_i(U_i) \subset [0, 1]^{\tilde{k}}$ .

**Step 3: Approximate Lipschitz mapping  $f \circ \phi_i^{-1}$  by  $\hat{f}_i$**

To approximate  $f \circ \phi_i^{-1} : [0, 1]^{\tilde{k}} \rightarrow \mathbb{R}$  via ReLU networks, we first caculate the Lipschitzness of it. For any pair  $x, y$  of  $\phi_i(U_i)$ , we have

$$\begin{aligned} |f \circ \phi_i^{-1}(x) - f \circ \phi_i^{-1}(y)| &\leq L \|\phi_i^{-1}(x) - \phi_i^{-1}(y)\|_\infty \\ &\leq L \|\phi_i^{-1}(x) - \phi_i^{-1}(y)\|_2 \\ &\leq \frac{2L}{1 - \delta} \sqrt{\frac{d}{\tilde{k}}} \|x - y\|_2 \\ &\leq \frac{2L\sqrt{d}}{1 - \delta} \|x - y\|_\infty. \end{aligned}$$

The first inequality is due to the Lipschitzness of function  $f$ . The second and last equality is the equivalence between  $\ell_2$  norm and  $\ell_\infty$  norm. The third inequality uses the isotropic property of the orthoprojector  $A_i$ . So  $f \circ \phi_i^{-1}$  is a  $\frac{2L\sqrt{d}}{1 - \delta}$  Lipschitz mapping from  $[0, 1]^{\tilde{k}}$  to  $\mathbb{R}$ .

By using Lemma C.2, there exists a ReLU network  $\hat{f}_i$  with at most

$$c_1 \left( \frac{c_2 \epsilon_1 (1 - \delta)}{2L\sqrt{d}} \right)^{-\tilde{k}} (\tilde{k}^2 + \tilde{k} \log \tilde{k} + \tilde{k} \log \frac{1}{\epsilon_1})$$

parameters such that for any  $x \in \phi_i(U_i)$ , we have the uniform error  $\epsilon_1$  as

$$|f \circ \phi_i^{-1}(x) - \hat{f}_i(x)| \leq \epsilon_1.$$

Notice that  $\phi_i$  is a linear mapping so that we can use a ReLU network with only one layer to represent it, which shows that we can approximate  $f$  efficiently in the local neighborhood  $U_i$ .

**Step 4: Approximate simple piecewise polynomial  $\rho_i \circ T_i^{-1}$  by  $\hat{\rho}_i$**

According to the poly-partitionable property of manifold  $M$  and Lemma B.2, there exists a ReLU network  $\hat{\rho}_i$  with at most  $O(k \log(k/\epsilon_2))$  parameters such that for any  $x \in T_i(U_i) \subset [0, 1]^k$ , we have the uniform error  $\epsilon_2$  as

$$|\rho_i \circ T_i^{-1}(x) - \hat{\rho}_i(x)| \leq \epsilon_2,$$

where  $T_i$  is composed by the tangent vectors of  $c_i$  and is scaled and translated to ensure  $T_i(U_i) \subset [0, 1]^k$ .

**Step 5: Determine the corresponding neighborhood for input**

Notice that  $\text{supp}(\rho_i) \subset U_i$  but  $\hat{\rho}_i$  may be non-zero for some point  $[0, 1]^k / T_i(U_i)$ , so we need to determine the corresponding chart for input  $x \in \mathcal{M}$  by ReLU networks. Inspired by Chen et al. (2019), we construct indicate approximator  $\hat{I}_\theta$  and  $\ell_2$  distance approximators  $\{\hat{d}_i^2\}_{i \in [N]}$  based on quadratic approximator in Lemma B.1 to approximate the neighborhood's indicator  $I\{x \in U_i\}$ , which relies upon the following identical equations

$$I\{x \in U_i\} = I\{\|x - c_i\|_2^2 < r^2\} = I\{(\cdot) < r^2\} \circ d_i^2(x),$$

where  $d_i^2(x)$  denotes the square of  $\ell_2$  distance between  $x$  and  $c_i$ . Then, if  $\hat{I}_\theta \approx I\{(\cdot) < r^2\}$  and  $\hat{d}_i^2 \approx d_i^2$ , we have  $\hat{I}_\theta \circ \hat{d}_i^2 \approx I\{(\cdot) < r^2\} \circ d_i^2 = I\{x \in U_i\}$ , which determines the corresponding chart approximately.

Assume that the uniform error of square distance approximator is  $\epsilon_q$  (i.e.  $|d_i^2 - \hat{d}_i^2| \leq \epsilon_q$  for any  $x \in [0, 1]^d$ ). In fact, functions computed by ReLU networks are piecewise linear but the indicator functions are not continuous, so we need to relax the indicator such that  $\hat{I}_\theta(x) = 1$  for  $x \leq r^2 + \epsilon_q - \theta$ ,  $\hat{I}_\theta(x) = 0$  for  $x \geq r^2 - \epsilon_q$  and  $\hat{I}_\theta$  is linear in  $(r^2 + \epsilon_q - \theta, r^2 - \epsilon_q)$ .

To correct the difference between indicator and its approximator, we will bound the value of function  $f$  such that the magnitude of  $f(x)$  is sufficient small when  $x$  is nearly on the boundary of  $U_i$ . Intuitively, for any  $y \in \partial(U_i)$ , we have

$$f\rho_i(y) = 0.$$

This is due to  $\text{supp}(\rho_i) \subset U_i$ , which implies that we only need estimate the upper bound of  $\|x - y\|_2$  for the Lipschitzness of  $f$  and smoothness of  $\rho_i$ , where  $x$  is nearly on  $\partial U_i$ . Indeed, we can prove that for any  $xU'_i := U_i/B_{\sqrt{r^2-\theta}}(c_i)$ , there exists  $y \in \partial U_i$  such that  $\|x - y\|_2 = O(\theta)$  (Chen et al., 2019, Lemma 3).

### Step 6: Estimate the total error

We combine three sub-network groups as

$$\hat{f} = \sum_{i=1}^N (\hat{f}_i \circ \phi_i) \hat{\times} (\hat{\rho}_i \circ T_i) \hat{\times} (\hat{I}_\theta \circ \hat{d}_i^2).$$

Next, we estimate the total error between  $f$  and  $\hat{f}$ . For any  $x \in M$ , we use  $g_i$  to denote  $(\hat{f}_i \circ \phi_i) \hat{\times} (\hat{\rho}_i \circ T_i)$ ,  $I_i$  to denote  $I\{x \in U_i\}$  and  $\hat{I}_i$  to denote  $\hat{I}_\theta \circ \hat{d}_i^2$ , then we have

$$\begin{aligned} |f(x) - \hat{f}(x)| &= \left| \sum_{i=1}^N f\rho_i - \sum_{i=1}^N g_i \hat{\times} \hat{I}_i \right| \\ &\leq \left| \sum_{i=1}^N f\rho_i - g_i I_i \right| + \left| \sum_{i=1}^N g_i \times I_i - g_i \hat{\times} \hat{I}_i \right| \\ &= \left| \sum_{i:x \in U_i} f\rho_i - (\hat{f}_i \circ \phi_i) \hat{\times} (\hat{\rho}_i \circ T_i) \right| + \left| \sum_{i=1}^N g_i \times I_i - g_i \hat{\times} \hat{I}_i \right| \\ &\leq \left| \sum_{i:x \in U_i} ((f \circ \phi_i^{-1} - \hat{f}_i) \circ \phi_i) \rho_i \right| + \left| \sum_{i:x \in U_i} (\hat{f}_i \circ \phi_i) \times \rho_i - (\hat{f}_i \circ \phi_i) \hat{\times} (\hat{\rho}_i \circ T_i) \right| \\ &\quad + \left| \sum_{i=1}^N g_i \times I_i - g_i \hat{\times} \hat{I}_i \right|. \end{aligned}$$

The second identical equation is due to  $\text{supp}(\rho_i) \subset U_i$ . Notice that  $\sum_{i \in [N]} \rho_i = 1$ , then the first term satisfies that

$$\left| \sum_{i:x \in U_i} ((f \circ \phi_i^{-1} - \hat{f}_i) \circ \phi_i) \rho_i \right| \leq \left( \sum_{i:x \in U_i} \rho_i \right) \max_{i:x \in U_i} \{|f \circ \phi_i^{-1} - \hat{f}_i|\} \leq \epsilon_1.$$

By the approximation of  $\hat{\times}$ , the second term satisfies that

$$\left| \sum_{i:x \in U_i} (\hat{f}_i \circ \phi_i) \times \rho_i - (\hat{f}_i \circ \phi_i) \hat{\times} (\hat{\rho}_i \circ T_i) \right| \lesssim \left| \sum_{i:x \in U_i} (\hat{f}_i \circ \phi_i) \times ((\rho_i \circ T_i^{-1} - \hat{\rho}_i) \circ T_i) \right| \leq c_f N \epsilon_2.$$

where  $c_f$  is the uniform upper bound the value of  $\{\hat{f}_i\}_{i \in [N]}$ . And the third term satisfies that

$$\left| \sum_{i=1}^N g_i \times I_i - g_i \hat{\times} \hat{I}_i \right| \lesssim \left| \sum_{i=1}^N g_i \times (I_i - \hat{I}_i) \right| \leq \sum_{i=1}^N \max_{x \in U'_i} |g_i| \lesssim \sum_{i=1}^N \max_{x \in U'_i} |f\rho_i| = O(N\theta).$$

Finally, we choose  $\epsilon_1 = O(\epsilon)$  and  $\epsilon_2 = \theta = O(\epsilon/N)$  to control the total error bounded by  $\epsilon$  and derive the upper bound for the size of network in Lemma E.1.

□



## E.2 Proof of Theorem 5.8

**Theorem E.4** (Restatement of Theorem 5.8). *Let  $\epsilon \in (0, 1)$  be a small constant. There exists a sequence  $\{N_k\}_{k \geq 1}$  that satisfies  $N_k = \Omega\left((2\epsilon\sqrt{d/k})^{-\frac{k}{2}}\right)$ , and a universal constant  $C_1 > 0$  such that the following holds: let  $\mathcal{M} \subset [0, 1]^d$  be a complete and compact  $k$ -dimensional Riemannian manifold with non-negative Ricci curvature, then there exists two  $2\epsilon$ -separated sets  $A, B \subset \mathcal{M}$  under  $\ell_\infty$  norm, such that for any  $\mu_0$ -balanced distribution  $P$  on the supporting set  $S = A \cup B$  and robust radius  $c \in (0, 1)$ , we have*

$$\inf \{ \mathcal{L}_P^{\infty, c\epsilon}(f) : f \in F_{N_k} \} \geq C_1 \mu_0.$$

*Proof.* Our proof relies on the following propositions.

**Lemma E.5.** (*Niyogi et al., 2008, Proposition 6.3*)

*Let  $\mathcal{M}$  be a sub-manifold of  $\mathbb{R}^d$  with condition number  $1/\tau$ . Let  $p$  and  $q$  be two points in  $\mathcal{M}$  such that  $\|x - y\|_2 = r$ . Then for all  $r \leq \tau/2$ , the geodesic distance  $d_{\mathcal{M}}(p, q)$  is bounded by*

$$d_{\mathcal{M}}(x, y) \leq \tau - \tau\sqrt{1 - 2r/\tau}.$$

By Lemma E.5, we know that  $d_{\mathcal{M}}(x, y) \leq \tau - \tau\sqrt{1 - 2r/\tau} \leq 2r$  when  $r \leq \tau/2$ .

**Lemma E.6.** (*Bishop, 1964, Bishop-Gromov Volume Comparison Theorem*) *Let  $\mathcal{M}$  is a complete Riemannian manifold with Ricci curvature  $\text{Ric} \geq (k - 1)l$ , and  $p \in \mathcal{M}$  is an arbitrary point. Then the function*

$$r \mapsto \frac{\text{vol}(B_{\mathcal{M}, r}(p))}{\text{vol}(B_r^l)}$$

*is a non-increasing function which tends to 1 as  $r$  goes to 0, where  $B_{\mathcal{M}, r}(p)$  is the  $\mathcal{M}$ 's geodesic ball of radius  $r$  and center  $p$ , and  $B_r^l$  is a geodesic ball of radius  $r$  in the space form  $\mathcal{M}_l^k$ . In particular,  $\text{vol}(B_{\mathcal{M}, r}(p)) \leq \text{vol}(B_r^l)$ .*

By Lemma E.6 and the non-negativeness of  $\mathcal{M}$ 's Ricci curvature, we know  $\text{vol}(B_{\mathcal{M}, r}(c)) \leq \text{vol}(B_r^0) = r^k V_k$ , where  $V_k$  denotes the volume of the unit ball in  $\mathbb{R}^k$ . Recall the relation between the covering number  $\mathcal{N}_{\mathcal{M}}(r)$  and the packing number  $\mathcal{P}_{\mathcal{M}}(r)$  on the manifold  $\mathcal{M}$ , then we have

$$\mathcal{P}_{\mathcal{M}}(r) \geq \mathcal{N}_{\mathcal{M}}(2r) \geq \frac{\text{vol}(\mathcal{M})}{(2r)^k V_k} = \Omega\left(\frac{\text{vol}(\mathcal{M}) k^{\frac{k}{2}}}{r^k}\right).$$

By choosing  $r = 2\epsilon\sqrt{d}$ , we know that there are at least  $\Omega\left((2\epsilon\sqrt{d/k})^{-k}\right)$  points on  $\mathcal{M}$  such that the  $\ell_\infty$  distance between each pair points of these is more than  $2\epsilon$ , where we use  $\mathcal{Q}$  to denote the set of these selected points. The remain of proof is similar to the latter half of proof for Theorem D.1.

Let  $S = \mathcal{Q}$  be the supporting set. Assume that for any partition  $A, B$  of  $S$  such that  $A \cup B = S$  and  $A \cap B = \emptyset$ , there exists a classifier  $f \in F_{N_k}$  that robustly classifies  $A$  and  $B$  with at least  $1 - \alpha$  accuracy. Next, we estimate the lower and upper bounds for the cardinal number of the vector set

$$R := \{(f(x))_{x \in \mathcal{Q}} | f \in F_{N_k}\}.$$

Let  $n$  denote  $|\mathcal{Q}|$ , then we have

$$R = \{(f(x_1), f(x_2), \dots, f(x_n)) | f \in F_{N_k}\},$$

where  $\mathcal{Q} = \{x_1, x_2, \dots, x_n\}$ .

On one hand, we know that for any  $u \in \{-1, 1\}^n$ , there exists a  $v \in R$  such that  $d_H(u, v) \leq \alpha n$ , where  $d_H(\cdot, \cdot)$  denotes the Hamming distance, then we have

$$|R| \geq \mathcal{N}(\{-1, 1\}^n, d_H, \alpha n) \geq \frac{2^n}{\sum_{i=0}^{\alpha n} \binom{n}{i}}.$$

On the other hand, by applying Lemma A.4, we have

$$\frac{2^n}{\sum_{i=1}^{\alpha n} \binom{n}{i}} \leq |R| \leq \Pi_{F_{N_k}}(n) \leq \sum_{j=0}^l \binom{n}{j}.$$

where  $l$  is the VC-dimension of  $F_{N_k}$ . In fact, we can derive  $l = \Omega(n)$  when  $\alpha$  is a small constant. Assume that  $l < n - 1$ , then we have  $\sum_{j=0}^l \binom{n}{j} \leq (en/l)^l$  and  $\sum_{i=1}^{\alpha n} \binom{n}{i} \leq (e/\alpha)^{\alpha n}$ , so

$$\frac{2^n}{(e/\alpha)^{\alpha n}} \leq |R| \leq (en/l)^l.$$

We define a function  $h(x)$  as  $h(x) = (e/x)^x$ , then we derive

$$2 \leq \left(\frac{e}{\alpha}\right)^\alpha \left(\frac{e}{l/n}\right)^{l/n} = h(\alpha)h(l/n).$$

When  $\alpha$  is sufficient small,  $l/n \geq C(\alpha)$  that is a constant only depending on  $\alpha$ , which implies  $l = \Omega(n)$ . Finally, by using Lemma A.3 and  $n = |\mathcal{Q}| = \Omega\left((2\epsilon\sqrt{d/k})^{-k}\right)$ , we know  $N_k = \Omega\left((2\epsilon\sqrt{d/k})^{-\frac{k}{2}}\right)$ . Combined with the definition of balanced distribution, we conclude the proof of Theorem E.4.  $\square$